

## Ensemble Technique in Decision Tree

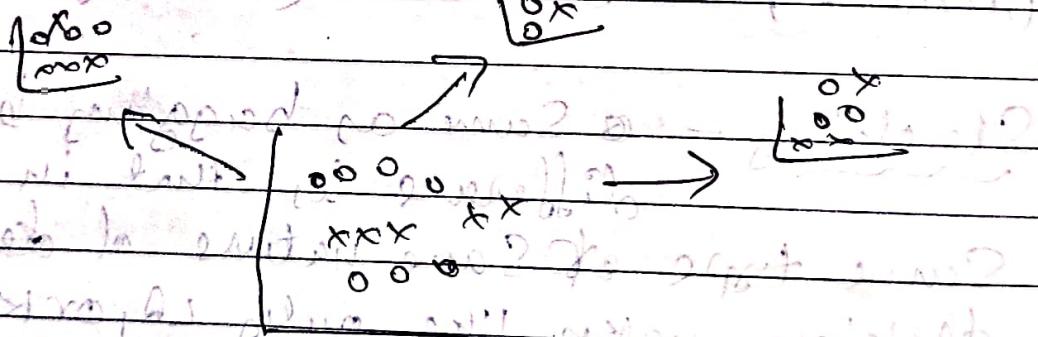
Date \_\_\_\_\_  
Page \_\_\_\_\_

- Ensemble Techniques → It is a technique where involve multiple decision maker. (Used to control the overfitting in Decision Tree)
- Bagging → In it multiple inside decision maker. Here every one will make decision and final decision will be there. (Same nature of decision maker)
- Boosting - Multiple decision maker but one will be a weak Scale decision maker, next one will better than previous, till continue. (Homogeneous Algo)
- Stacking - Same as bagging but the difference is that in bagging same type of same nature of decision maker like only LR, or KNN in bagging. But in Stacking all could be different one could be LR one could be KNN.

Note :- In Bagging homogeneous Algorithms are used whereas in Stacking heterogeneous Algorithms.

• Ensemble technique - one technique that  
combines multiple models to improve  
the accuracy of results by  
combining multiple model instead of using  
single model.

• Bootstrap Aggregation (Bagging) - It is a  
method which we divided our dataset into  
different, different sample with  
overlap or replacement.



• it decrease the Variance without  
increases of Bias. Because it is  
dividing into multiple dataset sample

• here data might overlap.

In Padding technique the data will not  
overlap.

# Random-forest

Date \_\_\_\_\_  
Page \_\_\_\_\_

# Bootstrap tech is also known as Bagging.

Bagging

SVM

KNN

Naive Bay's

Decision Tree

Random Forest

Ensemble Technique

AdaBoosting

Adaboost

Xg Boost

Gradient Boost

Stacking

## Random-Forest

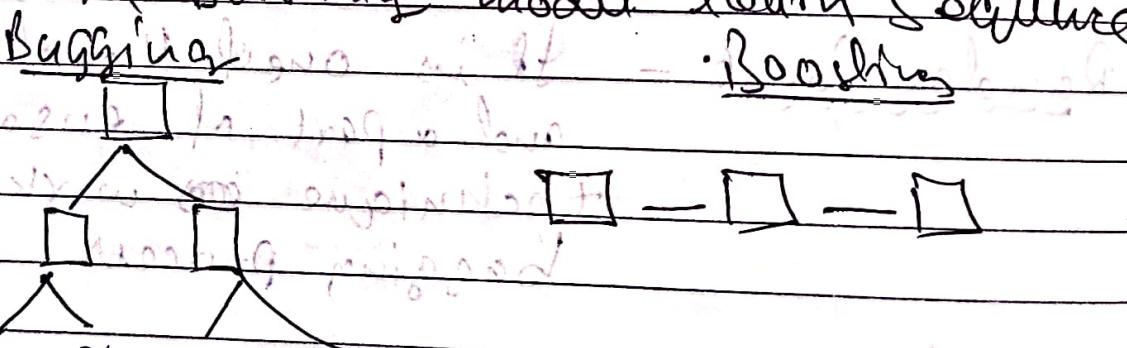
Random-Forest - It is one of the algorithm and a part of ensemble technique ~~to~~ work with bagging process.

Note → ~~as~~ we want to take ML model we always try to take a model which have low bias & low variance but to ~~get~~ get two simultaneous is very hard to hard and difficult because when bias is low then ~~the~~ variance is high & when variance is low then bias is high. So, To solve this problem we use Diversification, Bagging & Boosting.

Note → In Bagging we should always try to take a model who has low Bias & High Variance ( $L_2$  &  $H_V$ ) like (Decision Tree) (KNN) (Gully) but

In Boosting we should always try to take a model who has High Bias & Low Variance for ex (shallow decision tree) (depth should be very less), (decision stump)

Note — In Bagging model learn parallel & In Boosting model learn sequence



Note — Ensemble Technique is used to control the Overfitting in Decision Tree.

## Boosting

Date \_\_\_\_\_  
Page \_\_\_\_\_

$$\text{Boosting} = w_y = w_0 + \gamma \Delta F_w$$

$\gamma \rightarrow \text{city}$

$\Delta \rightarrow \text{letter}$

Boosting  $\rightarrow$  In Boosting we combine weak learners to make a overall strong learner.

## Ada. Boosting

Ex -	$f_1$	$f_2$	$f_3$	0/p   Weight   Update weight
-	-	-	Yes   $\frac{1}{7}$	0.058
-	-	-	No   $\frac{1}{7}$	0.058
-	-	-	Yes   $\frac{1}{7}$	0.058
incorrect record $\leftarrow$	-	-	$\uparrow$ -   $\frac{1}{7}$	0.349
-	-	-	-   $\frac{1}{7}$	0.058
-	-	-	-   $\frac{1}{7}$	0.058
-	-	-	-   $\frac{1}{7}$	0.058

$\therefore$  Here weight is created for Ada. Boosting

$$TE = \text{Total Error} = \frac{1}{7}$$

$$\text{Performance of Shimp (P.S)} = \frac{1}{2} \log_e \left( \frac{1-TE}{TE} \right)$$

$$PS \rightarrow \frac{1}{2} \log_e \left( \frac{1 - V_f}{V_f} \right) = 0.895$$

Now for New Sample —

for,

Correct Records

$$\Rightarrow \text{old weight} \times e^{-P_s}$$

$$= \frac{1}{2} \times e^{-0.895}$$

$$\Rightarrow 0.05$$

and

for incorrect record

and

we will ~~decide~~

in new sample

of weight

for, Incorrect Record

$$\Rightarrow \text{weight} \times e^{+15}$$

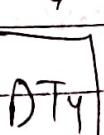
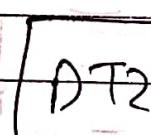
$$\Rightarrow \frac{1}{2} \times e^{+0.895}$$

$$\Rightarrow 0.349$$

Stump

Stump

Stump



Weak learners

Output

0

1

1

Majority Voting  $\Rightarrow$  Classification  $S_{avg} = 1$

Mean/Average  $\Rightarrow$  Regression  $L_{avg} = 3/4$

Q what is the white box & Black box

LR / Decision Tree

white Box

Ada Boost

Black box

# In Deep Learning, even Neural Network is used as a black box.

② Data package  $\rightarrow$  Dataset  $\xrightarrow{\text{Train}} \text{Model}$

Dataset  $\xrightarrow{\text{Test}}$

## Gradient - boosting - tree

- i) Gradient Boosting - It is a Success of Adaptive Boosting where the approach is little different from Ada Boost.

New Value = Old Value + Learning rate \* residual  
 where Old Value is average of outcome taken  
 Learning rate range from  $(0 - 0.00001)$

- We always try to keep residual as low as possible.

• Note → Gradient Boosting used for both Regression & Classification problems.  
 Here  $L(\alpha)$  is used to reduce overfitting.

- Prediction Residual in G.B is the kind of error

$$f(n) = h_0(n) + \alpha_1 h_1(n_1) + \alpha_2 h_2(n_2) + \dots + \alpha_n h_n(n_n)$$

base model       $\rightarrow DT_1$        $\rightarrow DT_2$

hypothesis      "      "      "      learning rate  $[0-1]$  (kind of hyperparameter)

- In boosting or in G.B we use  $\alpha$  (Alpha), learning rate for weak learner to strong learners.

$$f(n) = \text{Base Model} + \alpha_1(DT_1) + \alpha_2(DT_2) + \dots + \alpha_n(DT_n)$$

## XGBoost Algorithm -

- It improves the gradient boosting every further.
- The objective function of the XGBoost is that, it is a function of a loss function + a regularization.

Note - In regression any function or (loss func) will be MSE or MAE or others.

- In classification problem entropy or information gain for loss function.

$$\text{Obj}(\theta) = \sum_{i=1}^N l(y_i - \hat{g}_i) + \sum_{j=1}^J \alpha_j (f_j)$$

$\therefore$  this is loss function      this is a regularization

Note :- XGBoost doesn't use entropy or hessian Indices. Instead it utilises gradient (the error term) and hessian for creating the trees.

## XG Boost Classifier

(0.5 -)

Salinity	Enrichment	Approval	Residual
< 50K	B	0	-0.5
< 50K	C	1	0.5
> 50K	B	0	-0.5
> 50K	C	1	0.5
> 50K	N	1	0.5

i) construct Base model

ii) Construct Tree with root node

iii) Calculate Similarity Weight

$$\Rightarrow \sum (\text{Residual})^2$$

$$(ii) \rightarrow \sum P_{ab} (1 - P_{ab}) + \lambda$$

Probability

## XG Boost Regression -

<u>Exp</u>	<u>Crop</u>	<u>Salary</u>	<u>Res 1</u>
2	Yes	40K	-11
2.5	Yes	42K	-9
3	<del>Yes</del> No	52K	1
4	No	60K	9
4.5	Yes	62K	11

$\text{Res 1} = \text{Average of Salaries} - \text{Salary}$

$$\textcircled{(ii)} \text{ Similarity weight} = \frac{(\sum \text{Residual})^2}{\sum \text{per No. of residual} + 1}$$

∴ Here  $\lambda$  is a hyperparameter tuning.

## Bagging

Date \_\_\_\_\_  
Page \_\_\_\_\_

Code 1: `from sklearn.ensemble import BaggingClassifier`

`>>> bag_dt = BaggingClassifier(DecisionTreeClassifier(), n_estimators=10)`

`>>> bag_dt.fit(x-train, y-train)`  $\Rightarrow$  No. of tree

Note — Internally Random forest uses Decision Tree but it is also possible to use some other model or other algorithm as a base model for a particular bag.

It is not necessary to use always Random Forest for Bagging we can use any Base algorithm for Bagging.

Q: Is it possible to use Bagging without Decision Tree or Random forest.

A: Yes, it is possible, as a base model we can use any other algorithm.

Code 2: `>>> from sklearn.neighbors import KNeighborsClassifier`

`>>> bag_knn = BaggingClassifier(KNeighborsClassifier(n_neighbors=10))`

`>>> bag_knn.fit(x-train, y-train)`

```
>>> from sklearn.ensemble import RandomForestClassifier  
>>> rf = RandomForestClassifier(n_estimators=5)  
>>> rf.fit(x_train, y_train)  
>>> rf.estimators_  
>>> plt.figure(figsize=(10,10))  
>>> tree.plot_tree(rf.estimators_[0], filled=True)
```

newly added code to show who he Decision make

```
>>> plt.figure(figsize=(10,10))  
>>> for i in range(len(rf.estimators_)):   
>>> tree.plot_tree(rf.estimators_[i], filled=True)
```

Reaching point = Google AI paper with code

## (Boosting)

### Ada - Boosting

- Weak Learner  $\rightarrow$  If it is a ML model which have a accuracy just

- Decision Stump  $\rightarrow$  A decision tree which have max depth of 1 or having only one nodes

- +1 and -1  $\rightarrow$  Justified at 1 & 0 in classification

- Ada Boost is a Stage wise additive method.

Stage wise  $\rightarrow$  call weak learner. Sequence wise additive method  $\rightarrow$  In Ada-Boost if we apply add multiply weak learner.

$$\{ \alpha_i = \text{alpha}_i \}$$

$$\text{Ada-Boost} = \text{Sign}(\alpha_1 * h_1(n) + \alpha_2 * h_2(n) + \alpha_3 * h_3(n))$$

$$\text{Error rate} = 1 - \text{accuracy}$$

$$\alpha = \frac{1}{2} \log \left( \frac{1 - \text{error}}{\text{error}} \right)$$

but here error is to add the weight of  $\left(\frac{1}{n}\right)$  of wrong prediction.

$\alpha$  is learning rate.

1. **Boosting** - It is a ensemble modeling technique which improve the predictive power by converting a number of weak learners to strong learner.

2. The principal behind boosting algorithms is that we first build a model on the training dataset and then build a second model to rectify the error present in the first model.

3. The most common estimator used with AdaBoost is decision tree with one level which means decision tree with only 1 stump.



→ Decision tree with one stump / AdaBoost

4 In Ada Boost it builds a model and give equal weight to all the data point. If the assign higher weight to the wrong classifier. Now, all the points with higher weight are given more important in the next model. It will keep training model until and unless a lower error is received.

Model 1  $\rightarrow$  Model 2  $\rightarrow$  Model 3  $\rightarrow$  Model 4

weight 1  $\rightarrow$  weight 2  $\rightarrow$  weight 3  $\rightarrow$  weight 4

Formula of Weight =  $\frac{1}{\text{error}}$

Code [  $\gg$  from sklearn.ensemble import AdaBoostClassifier  
 $\gg$  clf = AdaBoostClassifier(random\_state=98)

$\gg$  clf.fit(X-train, y-train)

$\gg$  clf.score(X-train, y-train)

$\gg$  clf.score(X-test, y-test)

Now if we want to do Hyper parameter tuning then

>> from sklearn.ensemble import RandomForestClassifier

>> clf = AdaBoostClassifier(random\_state=96,  
base\_estimator=RandomForestClassifier(random\_state=101), n\_estimators=101,  
learning\_rate=0.01)

>> clf.fit(X\_train, y\_train)

∴ Now calculate Score of it.

i. Weight or TE (Total Error) =  $\frac{1}{N}$

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

i. New Sample weight = Old weight  $\times e^{\pm PS}$

(+) use for correct classified sample  
(-) use for wrong sample

## Gradient Boosting

- 1 Gradient Boost algorithm helps us to minimize bias errors of the model.
  - 2 Random: in ensemble technique the weak learners combine to make a strong model.
  - 3 Gradient Boosting is powerful enough to find every non-linear relationships between our model target and features.

Gradient Boosting — It is one of the variants of ensemble methods where we create multiple weak model and combine them to get better performance as a whole.

We always try to keep residual as low as possible.

Gradient boosting is typically used with decision tree (especially in CART)

formula  $\Rightarrow h_0(x) + \alpha_1 h_1(x_1) + \alpha_2 h_2(x_2) + \dots + \alpha_n h_n(x_n)$

Base model  $\Rightarrow$   $DT_1 \rightarrow DT_2 \rightarrow \dots$

or

$\Rightarrow$  Base model +  $\alpha_1 DT_1 + \alpha_2 DT_2 + \dots$

## XG Boosting

1) XG Boosting - It is designed for speed, ease of use, and performance on large datasets. It does not require optimization of the parameters or tuning.

2) XG Boost Features -

- i) XG Boost offers regularization which allows us to control overfitting by introducing L1/L2 penalties on the weights.
- ii) It also has a block structure for parallel learning.
- iii) It has a execution speed of model performance.

Q) How would you explain XG Boost in an interview?

A) XG Boost is a robust library that can help us improve our ML model accuracy. It is based on gradient boosting and can be used to fit any decision tree-based model.

## Stacking

Note :- Bagging & Boosting are homogeneous weak learner model for ensemble whereas Stacking often consider heterogeneous weak learners.

- In Stacking we are supposed to do a two level of a data division in general. We are supposed to divided original dataset into first one division and then we are supposed to take a subset of that data as a training and a testing data set.

Code 1 -   

```
>>> train, val, train, test, val-test = train_test_split  
(X, y, test_size=0.50,  
random_state=30)
```

After this will we create `x-train, y-train`

`>>> x-train, x-test, y-train, y-test = train_test_split  
(train, test, test_size=0.15)`

Now we will use model

Code 2 [  $\ggg$  knn = KNeighborsClassifier ()  
 $\ggg$  knn.fit (x\_train, y\_train) ]

Code 3 [  $\ggg$  svc = SVC () ]

$\ggg$  svc.fit (x\_train, y\_train) ]

Code 4 [  $\ggg$  prediction\_knn = knn.predict (val\_train) ]

$\ggg$  prediction\_svc = svc.predict (val\_train) ]

∴ We will do prediction from the first split data that is from Code no 1 with val\_train.

Code 5 Input 3:

Input 3  $\Rightarrow$  np.column\_stack ([ prediction\_knn, prediction\_svc ]) ]

↳ library np stack with two np arrays

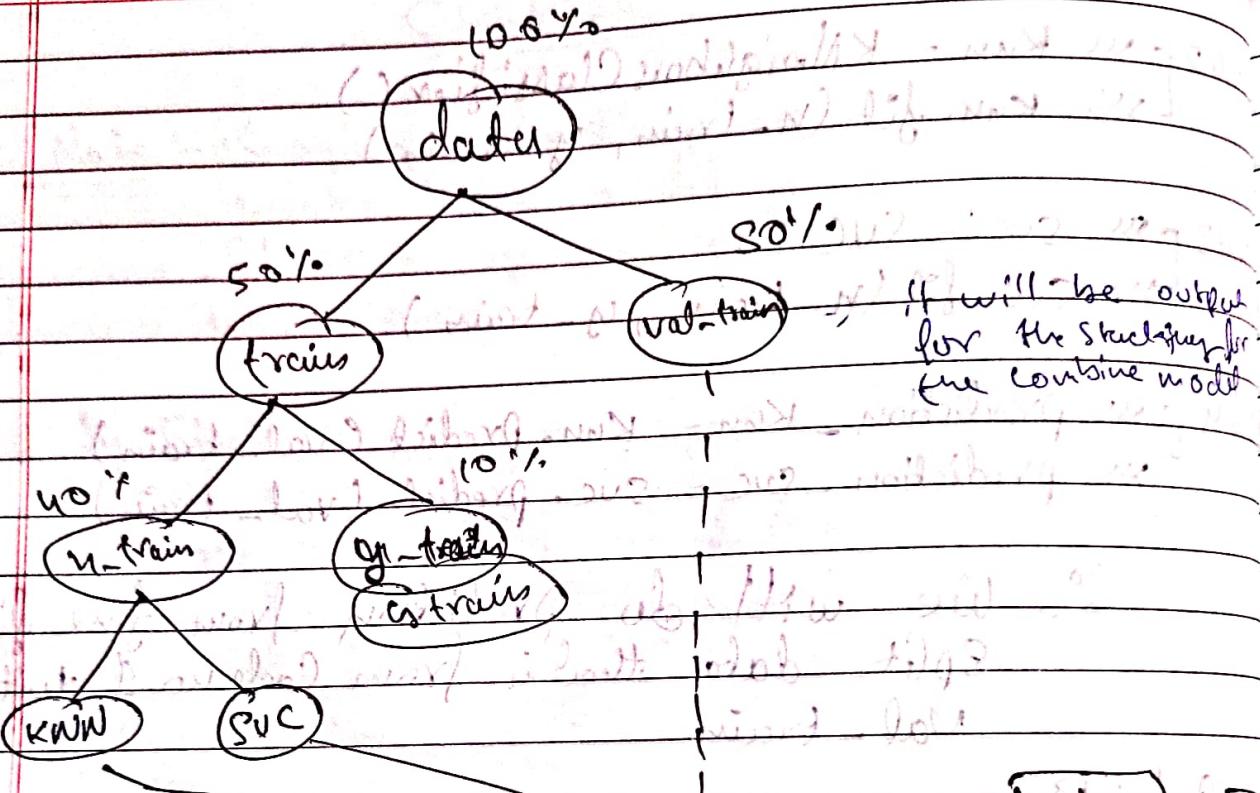
$\ggg$  output = val\_test ]

∴ Now we will create model 3.7 with the input 3 & output

Model 3.7 is linear together with gradient boosting

Want to use np.stack val if not

Note - We never pass a data to next Stack we are trying to pass a meta data / derived data of first Stack into a second Stack.



① we will use val-time data for prediction in KNN, SVC and after doing prediction of KNN & SVC we will combine their output as final output.

② After combining the output result we will use it for model 3 or next model.

code 6  
```rf = RandomForestClassifier()  
```rf.fit (input\_3, output)  
```rf.predict (input\_3)