

(Assignment-I) of Neha Verma (AIT). 8 - A

- Ans 1). a).
- Ans 2). c).
- Ans 3). d).
- Ans 4). b).
- Ans 5). c).
- Ans 6). c).
- Ans 7). b).
- Ans 8). d).
- Ans 9). c).
- Ans 10). c).

Q. Exp. overfitting in Decision tree using Depth as parameter?

A:- The reason behind overfitting in Decision tree due to the increasing no. of tree. We can solve it by pruning, hyperparameter tuning, max depth limit.

Then exp. how bagging & Random forest address this problem?

- In Bagging, we train on multiple models on different samples of Data.
- In Random forest, we usually combine multiple Decision tree in one random forest to get more accurate results.
- Due to this, the chances of overfitting will reduced.

Q:- Exp. Random forest working in details?

- A:-
- (A) Bootstrap Sampling
 - (B) Random feature selection
 - (C) Majority Voting

- A:- Basically, Random forest is that which use multiple decision tree to combine multiple Decision tree to get accurate result.
- Due to this, the chances of overfitting will reduce.
- Random forest uses bootstrap Sampling & Aggregation in which, we have some kind of Data, Divide random forest to Multiple Decision tree, the step before dividing tree is called Bootstrap & After averaging the value means we got 0.1, 0.9 anything that lead to aggregation Sampling.
- Majority Voting is used in classification in which, we assume two trees are giving us 1 or one tree giving us 0 val. So, there will be majority voting means we got 1 as our value.

Q:- A fraud Detection model produced the following results :-

	Predicted	not Predicted
Actual fraud	120	30
not fraud	50	800

Actual fraud 120 of out 300 is flagged as not fraud

not fraud 50 800

(a) for calculating Accuracy we use → Accuracy Score from sklearn.metrics import accuracy_score.
print ("Accuracy Score\n", accuracy_score(y_test, y_pred)).

(b) for calculating precision, recall & f-score. → we use confusion matrix.

from sklearn.metrics import accuracy_score, confusion_matrix.

TP = 120, FP = 30, FN = 30, TN = 800
 Total = 1000

$$\text{Accuracy} \rightarrow \frac{TP+TN}{\text{Total}} = \frac{120+800}{1000} = \frac{920}{1000} = 0.92$$

$$\text{Accuracy} \rightarrow 92\%$$

$$\text{Precision} \rightarrow \frac{TP}{TP+FP} = \frac{120}{120+30} = \frac{120}{150} = 0.705$$

$$\text{Precision} \rightarrow 70.5\%$$

$$\text{Recall} \rightarrow \frac{TP}{TP+FN} = \frac{120}{120+30} = 0.80$$

$$\text{Recall} \rightarrow 80\%$$

$$F_1\text{ score} \rightarrow \frac{2 \times \text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}} = \frac{2 \times 0.705 \times 0.80}{0.705 + 0.80} = \frac{1.128}{1.505}$$

$$f_1\text{ score} = 0.75$$

$$f_1\text{ score} \rightarrow 75\%$$

c) Model Acceptable Because

High Recall (80%) \rightarrow most fraud cases.

Better Accuracy we got (92%).

Q:- What happen if max. feature none?

A:- All features considered at every split..

Q:- Will this model overfit more or less than single Decision tree? Why?

(b) This model will overfit less than single Decision tree because we combine multiple Decision trees over there because of that we got good Accuracy. But in single Decision tree, ~~over~~ ^{one} model will definitely overfit because the max depth of tree is reached.

→ we can use Bootstrap Sampling.

→ Majority Voting.

→ Random feature selection

(c) Increase no. of tree if :-

Estimations → no. of tree (200)

Bias → slightly increases chance.

Variance → will be Decrease.

Because in Random forest we have :-

{ low Bias } required.

{ low variance }

Random Forest is better than (1) Most diff. (2) Top few words used

Random Forest is better than (1) Most diff. (2) Top few words used