

- 1) a) Large k with noisy data
- 2) a) They cannot handle missing value
and some from too large dataset they overfit
- 3) d) Improving feature selection
- 4) b) Features are ignored
- 5) a) Target variable is categorical
- 6) c) Sigmoid
- 7) b) Recall
- 8) c) Training time
- 9) a) To reduce bias
- 10) c) Logistic Regression

11) Decision Tree →

- Decision tree is used in supervised learning.
- It is a hierarchical structure.
- It includes root node, Parent node, Leaf node

Root node → First node that split the tree

Parent node → node used to split the tree

Leaf node → last node that predict the output

Now,

[Overfitting in Decision Tree using depth as a parameter]

Decision tree is a hierarchical structure. Some time it become overfit, by the increasing depth of the tree. because during the training time it learn the position of the data and give us high accuracy but from the testing time it can't handle the new data and become overfit.

means increasing the size of data having high chance of overfitting.

For solving that problem of overfitting we use Bagging and Random forest.

Bagging → Bagging is a technique of preventing model from overfitting by training multiple model from same data set and comparing their results.

Random forest → Random forest is a assembling technique it combine multiple decision trees and give us best result by combining multiple decision tree the size of data (depth of data) is divided into multiple decision

Trees and over model is solve the problem of overfitting.

→ Random forest →

Working →

- Random forest is assembling technique first it divide the dataset into different decision trees

- It use Random feature selection means the data set is divided randomly b/w the different decision tree.

- Every decision tree perform their operation like find the entropy and calculate the information gain

$$\text{Entropy} = -\sum p_i \log_2(p_i)$$

$$\text{IG} = \text{Entropy of parent} - \text{Weighted entropy of child}$$

- At last they combine each the decision trees and give the result by majority voting.

[Majority Voting → Count the result of each decision tree and give that result that are mostly given by the model]

~~→ In KNN, which situation will mas~~

	Predicted fraud	Predicted not fra
Actual fraud	120	30
Actual Not fraud	50	800

$$TP = 120$$

$$FP = 50$$

$$TN = 30$$

$$FN = 800$$

$$(a) \text{ Accuracy} = \frac{TP + TN}{TP + TN + FP + FN}$$

$$= \frac{120 + 30}{120 + 30 + 50 + 800}$$

$$= \frac{150}{900}$$

e) Yes

$$(b) \text{ Precision} = \frac{TP}{TP + FP}$$

$$= \frac{120}{120 + 50} = \frac{120}{170}$$

$$(c) \text{ Recall} = \frac{TP}{TP + FN} = \frac{120}{120 + 800} = \frac{120}{920}$$

$$(d) F_1 = \frac{\text{Precision}}{\text{Recall}}$$

$$= \frac{120}{170} \times \frac{920}{120} = \frac{920}{170}$$

- 13) a) Features are ignored
- b) Yes the model overfit because the fraction depth of the model = none
- c) Bias - high
variance - high.