day=13      Assessment -1

Answers

(1) a) large K with noisy data

(2) (c) Try overfit

(3) (c) Reduce variance

(4) (c) All features are considered at each split

(5) (c) Features are independent

(6) (c) ~~Sigmoid~~

(7) (d) $F_1$- score

(8) (d) Overfitting

9) a) to reduce bias

10) c) Logistic Regression

11) <u>Overfitting in Decision Trees :</u>

→ Overfitting is a condition when a model start learning a ~~data~~ training data rather than learning patterns, logic and relationship b/w target & features. and can't perform well on unseen data.

→ Decision trees are unstable learners because in most cases and depend upon data they prone to overfitting means low bias but high variance

→ generally, ~~depth~~ max-depth is a hyperparameter in decision trees which decides that how deep the splits are gone, if we select the high value for it the model can be prone to overfitting because all the focus will which biased to splitting maximum as possible and which cause cause overtraining on training data cause overfitting

→ solution, use bagging technique → Random Forest
It solves the problem of overfitting (low bias,
high variance) to (low bias, low variance)
by training n-estimators → no. of decision trees
parally on n-samples of data by using
technique called as bootstraping which
~~split~~ send samples of the original training data
by replacing rows and features depend on us.

→ and n-estimators learns on diff samples
of ~~training~~ data and produce outputs

→ and we use techniques like voting for classification
and average for numerical outputs of diff tree
and aggregate their result and produce
best model from no. of weak learners.
that's how bagging resolve the decision tree
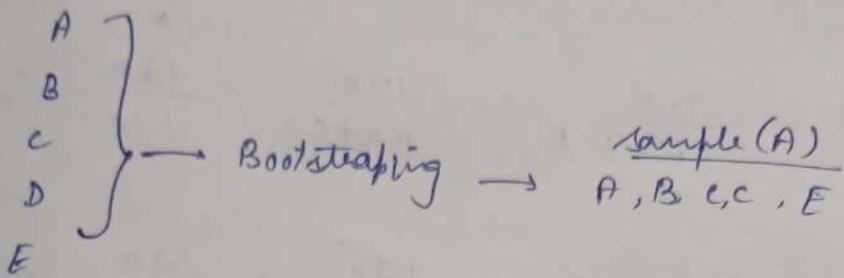problem.

<u>Ans:12) Random Forest Working:</u>

Random Forest is ~~a~~ Ensemble learning bagging
technique, which train no. of decision
tree on samples of data called weak
learners and ~~using~~ aggregate their outputs
to produce best model.

→ <u>Bootstrap sampling:</u> It is a sampling
data technique used by the random
forest to train no. of decision trees on
diff samples of same training data.
it uses row sampling and column sampling
to produce a data in a way that not
same data will repeat for all decision trees

like with row replacement → some rows can be
repeat in single sample
similarly around column replacement

eg: Row id

A
B
C } → Bootstrapping → sample (A)
D                        A, B, C, C, E
E

some rows can be repeat or some can never
be included ≈ 68% rows are selected and
32% can be item or for unseen data generally.
or without replacement → If a row is included
in a bag it will considered out of bag means
can't repeat now.

→ Random Feature selection: similarly some features
are included or can't be included for
different samples, features are selected randomly

→ Majority Voting: It is technique when we have
classification problem output select that one
for which most of the models produce
maximum same output                    Max Voting
                        output
eg:  Decision Tree 1 ————→ class 0
     DT 2           ————→ class 1  ✓  1 = 4
     DT 3           ————→ class 1  ✓  0 = 2
     DT 4           ————→ class 1  ✓
     DT 5           ————→ class 0
     DT 6           ————→ class 1 ✓

∴ Best Model will be output → class 1 selected
   as max weak learns produce class 1

Ans → 13

(a) Accuracy : $\dfrac{T.P + T.N}{T.P + T.N + F.N + F.P} = \dfrac{120 + 800}{120 + 800 + 50 + 30}$

$= \dfrac{920}{1000} = 0.92$

b) Precision : $\dfrac{T.P}{T.P + F.P} = \dfrac{120}{120 + 50} = \dfrac{120}{170} = 0.70$

c) Recall : $\dfrac{T.P}{T.P + F.N} = \dfrac{120}{120 + 30} = \dfrac{120}{150} = 0.8$

(d) $F_1$-score → $\dfrac{2 * \text{Precision} * \text{Recall}}{\text{Precision} + \text{Recall}} = \dfrac{2 * 0.70 * 0.8}{0.70 + 0.8}$

$= \dfrac{1.12}{1.5} = 0.74$

Actually these accuracy about 92%, but due to
confusion matrix actually 150 are actually
fraud out of 120 are positively classified as
fraud by model, 30 are mis classified,

So model can be considered but can be
cimproved for fraud detection.

Ans → 14

(i) if max_feature = None. → all features are
considered at each split

(ii) model will overfit more

(iii) If we choose n-estimators = 200, the
model will perform well on training data
and reduce overfitting while more

it will accept low variance and low bias → bias, variance trade-off.