

MACHINE LEARNING

1 In Q1 to Q7, only one option is correct, Choose the correct option:

1. The value of correlation coefficient will always be:
 - A) between 0 and 1
 - B) greater than -1
 - C) between -1 and 1**
 - D) between 0 and -1
2. Which of the following cannot be used for dimensionality reduction?
 - A) Lasso Regularisation
 - B) PCA
 - C) Recursive feature elimination**
 - D) Ridge Regularisation
3. Which of the following is not a kernel in Support Vector Machines?
 - A) linear
 - B) Radial Basis Function
 - C) hyperplane**
 - D) polynomial
4. Amongst the following, which one is least suitable for a dataset having non-linear decision boundaries?
 - A) **Logistic Regression**
 - B) Naïve Bayes Classifier
 - C) Decision Tree Classifier
 - D) Support Vector Classifier
5. In a Linear Regression problem, 'X' is independent variable and 'Y' is dependent variable, where 'X' represents weight in pounds. If you convert the unit of 'X' to kilograms, then new coefficient of 'X' will be?

(1 kilogram = 2.205 pounds)

 - A) $2.205 \times \text{old coefficient of 'X'}$
 - B) same as old coefficient of 'X'
 - C) old coefficient of 'X' $\div 2.205$**
 - D) Cannot be determined
6. As we increase the number of estimators in ADABOOST Classifier, what happens to the accuracy of the model?
 - A) remains same
 - B) increases**
 - C) decreases
 - D) none of the above
7. Which of the following is not an advantage of using random forest instead of decision trees?
 - A) Random Forests reduce overfitting
 - B) Random Forests explains more variance in data then decision trees
 - C) Random Forests are easy to interpret**
 - D) Random Forests provide a reliable feature importance estimate

In Q8 to Q10, more than one options are correct, Choose all the correct options:

8. Which of the following are correct about Principal Components?
 - A) Principal Components are calculated using supervised learning techniques
 - B) Principal Components are calculated using unsupervised learning techniques**
 - C) Principal Components are linear combinations of Linear Variables.**
 - D) All of the above
9. Which of the following are applications of clustering?
 - A) Identifying developed, developing and under-developed countries on the basis of factors like GDP, poverty index, employment rate, population and living index**
 - B) Identifying loan defaulters in a bank on the basis of previous years' data of loan accounts.**
 - C) Identifying spam or ham emails
 - D) Identifying different segments of disease based on BMI, blood pressure, cholesterol, blood sugar levels.**
10. Which of the following is(are) hyper parameters of a decision tree?
 - A) **max_depth**
 - B) max_features**
 - C) **n_estimators**
 - D) min_samples_leaf**

MACHINE LEARNING

Q10 to Q15 are subjective answer type questions, Answer them briefly.

11. What are outliers? Explain the Inter Quartile Range (IQR) method for outlier detection.

- An observation which differs from an overall pattern on a sample dataset is called an outlier.
- The outliers may suggest experimental errors, variability in a measurement, or an anomaly.
- Outliers affect mean and standard deviation of the dataset. These may statistically give erroneous results. Most machine learning algorithms do not work well in the presence of outlier.
- Thus, it is desirable to detect and remove outliers.
- Outliers are highly useful in anomaly detection like fraud detection where the fraud transactions are very different from normal transactions.

IQR is used to measure variability by dividing a data set into quartiles. The data is sorted in ascending order and split into 4 equal parts. Q1, Q2, Q3 called first, second and third quartiles are the values which separate the 4 equal parts.

Q1 represents the 25th percentile of the data.

Q2 represents the 50th percentile of the data.

Q3 represents the 75th percentile of the data.

If a dataset has $2n$ / $2n+1$ data points, then

Q1 = median of the dataset.

Q2 = median of n smallest data points.

Q3 = median of n highest data points.

IQR is the range between the first and the third quartiles namely Q1 and Q3: **$IQR = Q3 - Q1$.**

- **The data points which fall below $Q1 - 1.5 IQR$ or above $Q3 + 1.5 IQR$ are outliers.**

Find the lower and upper limits as $Q1 - 1.5 IQR$ and $Q3 + 1.5 IQR$

- **Data points greater than the upper limit or less than the lower limit are outliers**

12. What is the primary difference between bagging and boosting algorithms?

- Bagging is used when the goal is to reduce the variance of a decision tree classifier. Here the objective is to create several subsets of data from training sample chosen randomly with replacement. Each collection of subset data is used to train their decision trees. As a result, we get an ensemble of different models. Average of all the predictions from different trees are used which is more robust than a single decision tree classifier.
- Boosting is used to create a collection of predictors. In this technique, learners are learned sequentially with early learners fitting simple models to the data and then analysing data for errors. Consecutive trees (random sample) are fit and at every step, the goal is to improve the accuracy from the prior tree. When an input is misclassified by a hypothesis, its weight is increased so that next hypothesis is more likely to classify it correctly. This process converts weak learners into better performing model.

MACHINE LEARNING

13. What is adjusted R^2 in linear regression. How is it calculated?

R Square is a basic matrix which tells you about that how much variance is been explained by the model. What happens in a multivariate linear regression is that if you keep on adding new variables, the R square value will always increase irrespective of the variable significance. What adjusted R square do is calculate R square from only those variables whose addition in the model which are significant. So always while doing a multivariate linear regression we should look at adjusted R square instead of R square.

Calculation:

n = Total observations

p = Independent variables

$$\text{adj.r.squared} = 1 - (1 - \text{R.squared}) * ((n - 1)/(n - p - 1))$$

14. What is the difference between standardisation and normalisation?

- Normalization typically means rescales the values into a range of [0,1].
- Standardization typically means rescales data to have a mean of 0 and a standard deviation of 1 (unit variance).

15. What is cross-validation? Describe one advantage and one disadvantage of using cross-validation.

Cross-validation is a technique for evaluating ML models by training several ML models on subsets of the available input data and evaluating them on the complementary subset of the data. Use cross-validation to detect overfitting, ie, failing to generalize a pattern.

Advantages of Cross Validation

1. Reduces Overfitting: In Cross Validation, we split the dataset into multiple folds and train the algorithm on different folds. This prevents our model from overfitting the training dataset. So, in this way, the model attains the generalization capabilities which is a good sign of a robust algorithm.

Disadvantages of Cross Validation

1. Needs Expensive Computation: Cross Validation is computationally very expensive in terms of processing power required.