

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
Ans: C
2. Which of the following cannot be used for dimensionality reduction?
A) Lasso Regularisation B) PCA
C) Recursive feature elimination D) Ridge Regularisation
Ans: C
3. Which of the following is not a kernel in Support Vector Machines?
A) **linear** B) Radial Basis Function
C) hyperplane D) polynomial
Ans: A
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
Ans. A
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
Ans: A
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
Ans: C
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
Ans: A

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

MACHINE LEARNING

C) Principal Components are linear combinations of Linear Variables.

D) All of the above

Ans: B,C

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.

Ans: B,C

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

Ans: B,D

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.

Outliers:

The outliers may suggest experimental errors, variability in a measurement, or an anomaly. The age of a person may wrongly be recorded as 200 rather than 20 Years. Such an outlier should definitely be discarded from the dataset.

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.

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

Bagging is a way to decrease the variance in the prediction by generating additional data for training from dataset using combinations with repetitions to produce multi-sets of the original data. Boosting is an iterative technique which adjusts the weight of an observation based on the last classification.

MACHINE LEARNING

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

The adjusted R-squared is a modified version of R-squared that adjusts for the number of predictors in a regression model. It is calculated as:

$$\text{Adjusted } R^2 = 1 - [(1 - R^2) * (n - 1) / (n - k - 1)]$$

14. What is the difference between standardisation and normalisation?

Normalization is a scaling technique in which values are shifted and rescaled so that they end up ranging between 0 and 1. It is also known as Min-Max scaling.

Standardization is another scaling technique where the values are centered around the mean with a unit standard deviation. This means that the mean of the attribute becomes zero and the resultant distribution has a unit standard deviation.

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

Cross-Validation is a statistical method of evaluating and comparing learning algorithms by dividing data into two segments: one used to learn or train a model and the other used to validate the model.

Advantage:

More accurate estimate of out-of-sample accuracy.

The disadvantage of this method is that the training algorithm has to be rerun from scratch k times, which means it takes k times as much computation to make an evaluation. A variant of this method is to randomly divide the data into a test and training set k different times.
