

1. C
2. D
3. C
4. A
5. A
6. D
7. C
8. B,C
9. A,C,D
10. A,B,D

11. Outliers are values which are significantly far away from all the values in your data set which may occur because of difference in scale of measurement (like units) or due to human error.

Inter quartile range is the difference between the 25th percentile and 75th percentile of a data set.

Mathematically,

$IQR = Q3 - Q1$, where the Q3 is third quartile and Q1 is the first quartile.

According to the IQR method, any value beyond the range of $Q1 - 1.5 * IQR$ and $Q3 + 1.5 * IQR$ is considered an outlier.

12. Bagging and boosting both are ensemble methods. Bagging is a parallel method while boosting is a sequential method.

Bagging is a method to reduce the variation of the results thus reducing over-fitting. It is based on bootstrapping where we build various base models using bagging (taking random subsets of observations with replacement) and then take majority vote of all the base algorithms to create a strong model out of various weak models.

Boosting is based on giving more weightage to the misclassified observations in each next sequential model so as to give more emphasis to correctly classify the points which are hard to classify.

In boosting, subsets of data are fed into first base model which learns to predict and the misclassifications errors of first model are then fed into the next sequential model.

One more difference in boosting with respect to bagging is the samples drawn from data are not replaced.

13. Adjusted R-squared is an evaluation metrics used in regression models which is a modified version of R-squared adjusted to take into consideration the number of independent variables (predictors) used.

R-squared basically tells us about the amount of variance in our target variables explained by our independent features.

The mathematical formula for adjusted R-squared is as follows:

$$\text{Adj R-squared} = 1 - \{(n-1)/(n-k-1) * (RSS/TSS)\}$$

where n is total number of observations,

k is the number of features,

RSS is squared sum of differences between actual observed values and predicted values and

TSS is the squared sum of differences between actual observed values and mean of all the values.

14. Standardisation and Normalisation are techniques of data processing used to scale down the data to a same range.

The difference between standardisation and normalisation can be understood by looking at the mathematical implementation of both the techniques and the results obtained.

Standardisation means scaling down the data in such a way that its mean becomes zero and standard deviation becomes 1.

$$x_dash == (x_i - \text{mean}(X)) / \text{std}(X)$$

Normalisation refers to scaling down the data in such a way as to bring the data in the range of [0,1]

$$x_dash = x_i - \min(X) / \max(X) - \min(X)$$

15. Cross-validation is a technique of evaluating the algorithm by training it on multiple subsets of data thus leveraging the power of amount of data to result in better predicting models. We basically split the data into train, validation and test sets so that algorithm can be trained, validated and then tested on entirely unseen data.

In K-fold cross validation, we create equal folds of our train data and use (k-1) folds as our training data and 1 fold as the validation data in each of the k-iterations.

Advantages of Cross-validation:

Reduces over-fitting

Validates the performance of model by training it on several folds of data

Balances the classes in our data in case we are working with Imbalanced data set

Disadvantages of Cross-validation:

Training time increases by a factor of K as now we'll have to train the model K times using different subsets of data

Computational expense increases.