# Worksheet-Set 4: Machine learning Assignment

1. c
2. c
3. c
4. a
5. c
6. b
7. c
8. b,c
9. a,b,d
10. a,b,d

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

Outliers are data points that are significantly different from the majority of the data. They can be caused by errors in data collection, measurement, or entry, or they may represent genuinely unusual observations. Outliers can have a significant impact on statistical analyses, as they can distort the results and lead to incorrect conclusions. Therefore, it is important to identify and handle outliers appropriately.

One method for detecting outliers is the Inter Quartile Range (IQR) method. This method is based on the fact that most of the data in a normal distribution is contained within the range defined by the first quartile (Q1) and the third quartile (Q3). The IQR is calculated as the difference between Q3 and Q1.

To use the IQR method for outlier detection, we follow these steps:

1. Calculate Q1 and Q3. Calculate the IQR by subtracting Q1 from Q3.
2. Identify any points that fall outside of the range defined by Q1 minus 1.5 times the IQR and Q3 plus 1.5 times the IQR. These points are considered outliers.

For example, suppose that the Q1 and Q3 for a data set are 25 and 75, respectively. The IQR would be 75 - 25 = 50. Any points that fall outside of the range defined by 25 - (1.5 * 50) = -12.5 and 75 + (1.5 * 50) = 112.5 would be considered outliers.

The IQR method is only appropriate for identifying outliers in a distribution that is roughly symmetrical and bell-shaped. If the distribution is not symmetrical, other methods may be more appropriate for outlier detection.

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

Bagging and boosting are two ensemble learning methods that can be used to improve the performance of machine learning algorithms. Ensemble learning involves training a group of models and combining their predictions to make a more accurate overall prediction.

The primary difference between bagging and boosting algorithms is the way that they train and combine the individual models in the ensemble.

Bagging algorithms, also known as bootstrapped averaging algorithms, train a group of models independently on different subsets of the training data. The individual models are then combined by taking the mean or the majority vote of their predictions. Bagging algorithms are used to reduce the variance of a model, which can help to improve the stability and generalization of the model. Examples of bagging algorithms include random forests and bootstrapped decision trees.

Boosting algorithms, on the other hand, train a group of models sequentially, with each model attempting to correct the mistakes of the previous model. The individual models are combined by taking a weighted sum of their predictions. Boosting algorithms are used to reduce the bias of a model, which can help to improve the accuracy of the model. Examples of boosting algorithms include AdaBoost and XGBoost.

In general, bagging algorithms are more robust and easier to implement, but boosting algorithms can often achieve better performance on complex tasks.

## 13. What is adjusted R2 in linear regression. How is it calculated?

In linear regression, the coefficient of determination (R2) is a measure of how well the model fits the data. It is calculated as the ratio of the variance of the model's predictions to the variance of the actual data. Adjusted R2 is an adjusted version of R2 that takes into account the number of variables in the model. It is calculated as:

$$\text{Adjusted R2} = 1 - (1 - R2) * (n - 1) / (n - p - 1)$$

where n is the number of observations in the data set and p is the number of variables in the model.

Adjusted R2 is used to compare models with different numbers of variables and to determine how many variables are needed to achieve a good fit. It is generally preferred to R2 because it adjusts for the number of variables in the model, which can help to prevent overfitting.

To calculate adjusted R2, we will need to know the value of R2 for the model and the number of observations and variables in the data set. We can then use the formula above to calculate the adjusted R2.

## 14. What is the difference between standardization and Normalization?

Standardization and normalization are two techniques that are commonly used to scale numeric variables so that they can be compared or combined. Standardization involves scaling the variables so that they have a mean of 0 and a standard deviation of 1. This is done by subtracting the mean from each value and dividing the result by the standard

deviation. Standardization is often used when the variables are approximately normally distributed and there are no outliers.

Normalization, on the other hand, scales the variables so that they have a minimum value of 0 and a maximum value of 1. This is done by subtracting the minimum value from each value and dividing the result by the range (maximum value - minimum value). Normalization is often used when the variables are not normally distributed or there are outliers.

In general, standardization is more commonly used than normalization, but both techniques can be useful depending on the characteristics of the data and the needs of the analysis.

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

Cross-validation is a resampling procedure used to evaluate the performance of machine learning models. It involves dividing the data into a training set and a test set, training the model on the training set, and evaluating the model on the test set. This process is repeated a number of times, with different portions of the data being used as the test set each time. The results from each iteration are then averaged to estimate the overall performance of the model.

One advantage of cross-validation is that it allows you to get a more accurate estimate of the model's performance, as it uses more of the data for training and evaluation. This is particularly useful when the data set is small, as it can provide a more reliable estimate of the model's performance than a single train/test split.

One disadvantage of cross-validation is that it can be time-consuming, particularly for complex models that take a long time to train. In addition, cross-validation can be sensitive to the specific data splits used, so the results may vary depending on the way that the data is partitioned.