# MACHINE LEARNING
## WORKSHEET – 1

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

11. Definition of outliers

An outlier is an observation that lies an abnormal distance from other values in a random sample from a population. In a sense, this definition leaves it up to the analyst (or a consensus process) to decide what will be considered abnormal. Before abnormal observations can be singled out, it is necessary to characterize normal observations.

**Significance of outliers:**
- Outliers badly 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. So 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.

**What is Interquartile Range IQR?**

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. **Differences Between Bagging and Boosting –**

| S.NO | BAGGING | BOOSTING |
|------|---------|----------|
| 1. | Simplest way of combining predictions that belong to the same type. | A way of combining predictions that belong to the different types. |
| 2. | Aim to decrease variance, not bias. | Aim to decrease bias, not variance. |

| 3. | Each model receives equal weight. | Models are weighted according to their performance. |
|---|---|---|
| 4. | Each model is built independently. | New models are influenced by performance of previously built models. |
| 5. | Different training data subsets are randomly drawn with replacement from the entire training dataset. | Every new subsets contains the elements that were misclassified by previous models. |
| 6. | Bagging tries to solve over-fitting problem. | Boosting tries to reduce bias. |
| 7. | If the classifier is unstable (high variance), then apply bagging. | If the classifier is stable and simple (high bias) the apply boosting. |
| 8. | Random forest. | Gradient boosting. |

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

Adjusted R-squared  is used to compare the goodness-of-fit for       regression models that contain differing numbers of independent variables.       Logistic regression with binary and multinomial outcomes is commonly       used, and it is also searched for an interpretable measure of the strength of a particular logistic model.
    Let's have a look at the formula for adjusted R-squared to better understand its working.

Adjusted R-squared = {1 - [(1-R2)(n-1)/(n-k-1)]}

Here,
n represents the number of data points in our dataset
k represents the number of independent variables, and
R represents the R-squared values determined by the model

So, if R-squared does not increase significantly on the addition of a new independent variable, then the value of Adjusted R-squared will actually decrease.

On the other hand, if on adding the new independent variable we see a significant increase in R-squared value, then the Adjusted R-squared value will also increase. it is better to use Adjusted R-squared when there are multiple variables in the regression model.

14.

**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.**

Here's the formula for normalization:

$$X' = \frac{X - X_{min}}{X_{max} - X_{min}}$$

Here, Xmax and Xmin are the maximum and the minimum values of the feature respectively.

- When the value of X is the minimum value in the column, the numerator will be 0, and hence X' is 0
- On the other hand, when the value of X is the maximum value in the column, the numerator is equal to the denominator and thus the value of X' is 1
- If the value of X is between the minimum and the maximum value, then the value of X' is between 0 and 1

**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.**

- Here's the formula for standardization:

$$X' = \frac{X - \mu}{\sigma}$$

- $\mu$ is the mean of the feature values and $\sigma$ is the standard deviation of the feature values. Note that in this case, the values are not restricted to a particular range.

15.

Cross-validation is a technique in which we train our model using the subset of the data-set and then evaluate using the complementary subset of the data-set.

The three steps involved in cross-validation are as follows :

1. Reserve some portion of sample data-set.
2. Using the rest data-set train the model.
3. Test the model using the reserve portion of the data-set.

**Methods of Cross Validation**

**Validation**
In this method, we perform training on the 50% of the given data-set and rest 50% is used for the testing purpose. The major drawback of this method is that we perform training on the 50% of the dataset, it may possible that the remaining 50% of the data contains some important information which we are leaving while training our model i.e higher bias.

**LOOCV (Leave One Out Cross Validation)**
In this method, we perform training on the whole data-set but leaves only one data-point of the available data-set and then iterates for each data-point. It has some advantages as well as disadvantages also.
An advantage of using this method is that we make use of all data points and hence it is low bias.
The major drawback of this method is that it leads to higher variation in the testing model as we are testing against one data point. If the data point is an outlier it can lead

to higher variation. Another drawback is it takes a lot of execution time as it iterates over 'the number of data points' times.

**K-Fold Cross Validation**

In this method, we split the data-set into k number of subsets(known as folds) then we perform training on the all the subsets but leave one(k-1) subset for the evaluation of the trained model. In this method, we iterate k times with a different subset reserved for testing purpose each time.

*Note:*

It is always suggested that the value of k should be 10 as the lower value of k is takes towards validation and higher value of k leads to LOOCV method