

## **ASSIGNMENT - SUBJECTIVE QUESTIONS**

**1. From your analysis of the categorical variables from the dataset, what could you infer about their effect on the dependent variable?**

Answer : Univariate and bivariate analysis gave pretty good insights,

- Clear weather attracted more booking which seems obvious.
- Booking seemed to be almost equal either on working day or non-working day
- There is increase in rentals in 2019.
- There is high demand on working days.

**2. Why is it important to use drop\_first=True during dummy variable creation?**

Answer: drop\_first = True is important to use, as it helps in reducing the extra column created during dummy variable creation. Hence it reduces the correlations created among dummy variables.

**3. Looking at the pair-plot among the numerical variables, which one has the highest correlation with the target variable?**

Answer: 'atemp' and 'temp' has the highest correlation with the target variable

**4. How did you validate the assumptions of Linear Regression after building the model on the training set?**

Answer : Model with better R-squared and Adjusted R-squared value, Multicollinearity check, Residual check.

**5. Based on the final model, which are the top 3 features contributing significantly towards explaining the demand of the shared bikes?**

Answer : atemp, temp, humidity

## **GENERAL - SUBJECTIVE QUESTIONS**

### **1. Explain the linear regression algorithm in detail.**

Answer: Linear regression may be defined as the statistical model that analyses the linear relationship between a dependent variable with given set of independent variables. Linear relationship between variables means that when the value of one or more independent variables will change (increase or decrease), the value of dependent variable will also change accordingly (increase or decrease).

Mathematically the relationship can be represented with the help of following equation –  $Y = mX + c$

Here,

- Y is the dependent variable we are trying to predict.
- X is the independent variable we are using to make predictions.
- m is the slope of the regression line which represents the effect X has on Y
- c is a constant, known as the Y-intercept. If  $X = 0$ , Y would be equal to c

### **2. Explain the Anscombe's quartet in detail.**

Answer: Anscombe's Quartet was developed by statistician Francis Anscombe. It comprises four datasets, each containing eleven (x, y) pairs. The essential thing to note about these datasets is that they share the same descriptive statistics.

### **3. What is Pearson's R?**

Answer: Pearson's r is a numerical summary of the strength of the linear association between the variables. If the variables tend to go up and down together, the correlation coefficient will be positive. If the variables tend to go up and down in opposition with low values of one variable associated with high values of the other, the correlation coefficient will be negative.

### **4. What is scaling? Why is scaling performed? What is the difference between normalized scaling and standardized scaling?**

Answer: Feature Scaling is a technique to standardize the independent features present in the data in a fixed range. It is performed during the data pre-processing to handle highly varying magnitudes or values or units. If feature scaling is not done, then a machine learning algorithm tends to weigh greater values, higher and consider smaller values as the lower values, regardless of the unit of the values.

### **5. You might have observed that sometimes the value of VIF is infinite. Why does this happen?**

Answer: If there is perfect correlation, then  $VIF = \text{infinity}$ . A large value of VIF indicates that there is a correlation between the variables. If the VIF is 4, this means that the variance of the model coefficient is inflated by a factor of 4 due to the presence of multicollinearity.

**6. What is a Q-Q plot? Explain the use and importance of a Q-Q plot in linear regression.**

Answer: The quantile-quantile (q-q) plot is a graphical technique for determining if two data sets come from populations with a common distribution.

Use of Q-Q plot: A q-q plot is a plot of the quantiles of the first data set against the quantiles of the second dataset. By a quantile, we mean the fraction (or percent) of points below the given value. That is, the 0.3 (or 30%) quantile is the point at which 30% percent of the data fall below and 70% fall above that value. A 45-degree reference line is also plotted. If the two sets come from a population with the same distribution, the points should fall approximately along this reference line. The greater the departure from this reference line, the greater the evidence for the conclusion that the two data sets have come from populations with different distributions.