

## Assignment-based Subjective

**Questions 1. From your analysis of the categorical variables from the dataset, what could you infer about their effect on the dependent variable? (3 marks)**

Answer:

2019 attracted more number of booking from the previous year, which shows good progress in terms of business, Most of the bookings has been done during the month of may, june, july, aug, sep and oct. Trend increased starting of the year till mid of the year and then it started decreasing as we approached the end of year. Fall season seems to have attracted more booking. And, in each season the booking count has increased drastically from 2018 to 2019. Clear weather attracted more booking which seems obvious. Thu, Fir, Sat and Sun have more number of bookings as compared to the start of the week. When it's not holiday, booking seems to be less in number which seems reasonable as on holidays, people may want to spend time at home and enjoy with family.

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

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. Let's say we have 3 types of values in Categorical column and we want to create dummy variable for that column. If one variable is not furnished and semi\_furnished, then It is obvious unfurnished. So we do not need 3rd variable to identify the unfurnished.

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

'temp' variable 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?**

One of the ways to determine if this assumption is met or not is by creating a scatter plot x vs y. If the data points fall on a straight line in the graph, there is a linear relationship between the dependent and the independent variables, and the assumption holds.

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

- Holiday

- Temperature
- Humidity
- Season

## General Subjective Questions

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

Linear regression is a type of supervised machine learning algorithm that computes the linear relationship between a dependent variable and one or more independent features. When the number of the independent feature, is 1 then it is known as Univariate Linear regression, and in the case of more than one feature, it is known as multivariate linear regression. The goal of the algorithm is to find the best linear equation that can predict the value of the dependent variable based on the independent variables. The equation provides a straight line that represents the relationship between the dependent and independent variables. The slope of the line indicates how much the dependent variable changes for a unit change in the independent variable(s)

Linear regression is used in many different fields, including finance, economics, and psychology, to understand and predict the behavior of a particular variable. For example, in finance, linear regression might be used to understand the relationship between a company's stock price and its earnings or to predict the future value of a currency based on its past performance. One of the most important supervised learning tasks is regression. In regression set of records are present with X and Y values and these values are used to learn a function so if you want to predict Y from an unknown X this learned function can be used. In regression we have to find the value of Y, So, a function is required that predicts continuous Y in the case of regression given X as independent features. Here Y is called a dependent or target variable and X is called an independent variable also known as the predictor of Y. There are many types of functions or modules that can be used for regression. A linear function is the simplest type of function. Here, X may be a single feature or multiple features representing the problem. Linear regression performs the task to predict a dependent variable value (y) based on a given independent variable (x)). Hence, the name is Linear Regression. In the figure above, X (input) is the work experience and Y (output) is the salary of a person. The regression line is the best-fit line for our model

### 2) Explain the Anscombe's quartet in detail

Anscombe's quartet comprises a set of four dataset, having identical descriptive statistical properties in terms of means, variance, R-Squared, correlations, and linear regression lines but having different representations when we scatter plot on graph. The datasets were created by the statistician Francis Anscombe in 1973 to demonstrate the importance of visualizing data and to show that summary statistics alone can be misleading. The four datasets that make up Anscombe's quartet each include 11 x-y pairs of data. When plotted, each dataset seems to have a unique connection between x and y, with unique variability patterns and distinctive correlation strengths. Despite these variations, each dataset has the same summary statistics, such as the same x and y mean and variance, x and y correlation coefficient, and linear regression line.

Anscombe's quartet is used to illustrate the importance of exploratory data analysis and the drawbacks of depending only on summary statistics. It also emphasizes the importance of using data visualization to spot trends, outliers, and other crucial details that might not be obvious from summary statistics alone.

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

The Pearson correlation coefficient ( $r$ ) is the most widely used correlation coefficient and is known by many names: Pearson's  $r$ , Bivariate correlation, Pearson product-moment correlation coefficient (PPMCC), The correlation coefficient. The Pearson correlation coefficient is a descriptive statistic, meaning that it summarizes the characteristics of a dataset. Specifically, it describes the strength and direction of the linear relationship between two quantitative variables.

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

Normalization is a data preprocessing technique used to adjust the values of features in a dataset to a common scale. This is done to facilitate data analysis and modeling, and to reduce the impact of different scales on the accuracy of machine learning models. 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,  $X_{\max}$  and  $X_{\min}$  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. What is Standardization? Standardization is another scaling method 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: 
$$X' = \frac{X - \mu}{\sigma}$$
 where  $\mu$  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.

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

If there is perfect correlation, then  $VIF = \infty$ . 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.**

Q-Q plots are also known as Quantile-Quantile plots. As the name suggests, they plot the quantiles of a sample distribution against quantiles of a theoretical distribution. Doing this helps us

determine if a dataset follows any particular type of probability distribution like normal, uniform, exponential