# 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 - My analysis of categorical variables from the dataset on the dependent variables are:

- Weather is most optimal for bike renting bike rental is more in partly cloudy weather and during the fall season and then in summer
- Working and non-working days have almost the same median although spread the bigger for non-working days.
- Bike Rentals are more in the year 2019 compared to 2018
- People rent more on non-holidays compared to holidays, it has been also noticed that bike rental is more on Saturday, Wednesday, and Thursday

2. Why is it important to use **drop_first=True** during dummy variable creation? (2 mark)

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.
A variable with n levels can be represented by n -1 dummy variables. So, if we remove the first column then also, we represent the data. Id the value of variable from 2 to n is 0, it means that the value of $1^{st}$ variable is 1

3. Looking at the pair-plot among the numerical variables, which one has the highest correlation with the target variable? (1 mark)

Answer – temp had the highest correlation of 0.99

4. How did you validate the assumptions of Linear Regression after building the model on the

   training set? (3 marks)

Answer – By plotting the residuals distribution. It came out to be a normal distribution with a mean value of 0.

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

Answer - The top 3 features contribution towards explaining the demands of the shared bikes are :

1. holiday
2. yr
3. light rain

# General Subjective Questions

1. Explain the linear regression algorithm in detail.

Answer - It is a method of finding the best straight-line fitting to the given dataset With the help of linear regression algorithm we tries to find the best linear relationship between the independent and dependent variables.
It is a supervised machine learning algorithm that finds the best linear-fit relationship on the given dataset, between independent and dependent variables. It is done with the help of the Sum of Squared Residuals Method

2. Explain the Anscombe's quartet in detail.

Answer - Anscombe's quartet comprises four datasets that have nearly identical simple statistical properties yet appear very different when graphed. Each dataset consists of eleven (x,y) points. They were constructed in 1973 by the statistician Francis Anscombe to demonstrate both the importance of graphing data before analyzing it and the effect of outliers on statistical properties.

3. What is Pearson's R ?

Answer - the Pearson correlation coefficient (PCC), also referred to as Pearson's r, the Pearson product-moment correlation coefficient (PPMCC), or the bivariate correlation, is a measure of linear correlation between two sets of data. It is the covariance of two variables, divided by the product of their standard deviations; thus it is essentially a normalised measurement of the covariance, such that the result always has a value between −1 and 1.
The Pearson's correlation coefficient varies between -1 and +1 where:
• 	r = 1 means the data is perfectly linear with a positive slope ( i.e., both variables tend to change in the same direction)
• 	r = -1 means the data is perfectly linear with a negative slope ( i.e., both variables tend to change in different directions)
• 	r = 0 means there is no linear association
• 	r > 0 < 5 means there is a weak association
• 	r > 5 < 8 means there is a moderate association
• 	r > 8 means there is a strong association

4. What is scaling? Why is scaling performed? What is the difference between normalized scaling and standardized scaling? (3 marks)

Answer - *It is a step of data Pre-Processing which is applied to independent variables to normalize the data within a particular range. It also helps in speeding up the calculations in an algorithm.*

Normalization typically means rescales the values into a range of [0,1]. Standardization typically means rescales data to have a mean of 0 and a standard deviation of 1 (unit variance).

| S.NO. | Normalisation | Standardisation |
|---|---|---|
| 1. | Minimum and maximum value of features are used for scaling | Mean and standard deviation is used for scaling. |
| 2. | It is used when features are of different scales. | It is used when we want to ensure zero mean and unit standard deviation. |
| 3. | Scales values between [0, 1] or [-1, 1]. | It is not bounded to a certain range. |
| 4. | It is really affected by outliers. | It is much less affected by outliers. |
| 5. | Scikit-Learn provides a transformer called MinMaxScaler for Normalization. | Scikit-Learn provides a transformer called StandardScaler for standardization. |
| 6. | This transformation squishes the n-dimensional data into an n-dimensional unit hypercube. | It translates the data to the mean vector of original data to the origin and squishes or expands. |
| 7. | It is useful when we don't know about the distribution | It is useful when the feature distribution is Normal or Gaussian. |
| 8. | It is a often called as Scaling Normalization | It is a often called as Z-Score Normalization. |

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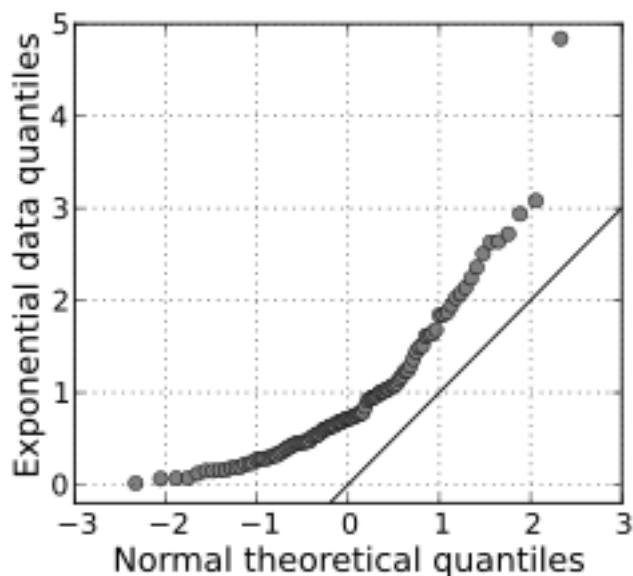
(3 marks)

Answer - If there is perfect correlation, then VIF = infinity. This shows a perfect correlation between two independent variables. In the case of perfect correlation, we get R2 =1, which lead to 1/(1-R2) infinity. To solve this problem, we need to drop one of the variables from the dataset which is causing this perfect multicollinearity.
An infinite VIF value indicates that the corresponding variable may be expressed exactly by a linear combination of other variables (which show an infinite VIF as well).

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

(3 marks)

Answer - Q-Q Plots (Quantile-Quantile plots) are plots of two quantiles against each other. A quantile is a fraction where certain values fall below that quantile. For example, the median is a quantile where 50% of the data fall below that point and 50% lie above it. The purpose of Q Q plots is to find out if two sets of data come from the same distribution. A 45 degree angle is plotted on the Q Q plot; if the two data sets come from a common distribution, the points will fall on that reference line.



If the two distributions being compared are similar, the points in the Q–Q plot will approximately lie on the line y = x. If the distributions are linearly related, the points in the Q–Q plot will approximately lie on a line, but not necessarily on the line y = x. Q–Q plots can also be used as a graphical means of estimating parameters in a location-scale family of distributions.