# 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:

- → The year box plots indicates that more bikes are rent during 2019.

- → The season box plots indicates that more bikes are rent during fall season.

- → The working day and holiday box plots indicate that more bikes are rent during normal working days than on weekends or holidays.

- → The month box plots indicates that more bikes are rent during September month.

- → The weekday box plots indicates that more bikes are rent during Saturday.

- → The weathersit box plots indicates that more bikes are rent during Clear, Few clouds, Partly cloudy weather.


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.
- → Example- We can create dummy variables on categorical variables season,yr, mnth, weekday, and weathersit.

3. Looking at the pair-plot among the numerical variables, which one has the highest correlation with the target variable?                                    (1 mark)
Answer:
- → By looking at the pair plot temp variable has the highest (0.63) correlation with target variable 'cnt(Count)'.

4. How did you validate the assumptions of Linear Regression after building the model on the training set?                                    (3 marks)
Answer:
- → From the histogram, we had see that the Residuals are normally distributed and mean is 0. Hence our assumption for Linear Regression is valid.

5. Based on the final model, which are the top 3 features contributing significantly towards explaining the demand of the shared bikes?                                    (2 marks)
Answer:
- → The Top 3 features contributing significantly towards the demands of share bikes are:
- → Light_rain_light_snow_thunderstroms(negative correlation).
- → yr_2019(Positive correlation).
- → Sept Month(Positive correlation).

# General Subjective Questions

1. Explain the linear regression algorithm in detail. (4 marks)

Answer:

Linear regression is one of the easiest and most popular Machine Learning algorithms. It is a statistical method that is used for predictive analysis. Linear regression makes predictions for continuous/real or numeric variables such as sales, salary, age, product price, etc.

Linear regression algorithm shows a linear relationship between a dependent (y) and one or more independent (y) variables, hence called as linear regression. Since linear regression shows the linear relationship, which means it finds how the value of the dependent variable is changing according to the value of the independent variable.

2. Explain the Anscombe's quartet in detail. (3 marks)
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.

3. What is Pearson's R? (3 marks)
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? (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.

Most of the times, collected data set contains features highly varying in magnitudes, units and range. If scaling is not done then algorithm only takes magnitude in account and not units hence incorrect modelling. To solve this issue, we have to do scaling to bring all the variables to the same level of magnitude.

It is important to note that scaling just affects the coefficients and none of the other parameters like t-statistic, F-statistic, p-values, R-squared, etc.

Normalization/Min-Max Scaling:

- It brings all of the data in the range of 0 and 1. sklearn.preprocessing.MinMaxScaler helps to implement normalization in python.

  MinMax Scaling: $x=x-min(x)/max(x)-min(x)$

Standardization Scaling:

- Standardization replaces the values by their Z scores. It brings all of the data into a standard normal distribution which has mean ($\mu$) zero and standard deviation one ($\sigma$).

  Standardization: $x= x-mean(x)/sd(x)$

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.