# 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?**

**Ans:** The demand of the bike is more in 2019 than 2018. In month 9, bike demand is higher compared to other months. Bike demand is low in spring season compared to other seasons. It is also seemed that, the bike demand is higher in weathersit1 (i.e., when sky is clear or few clouds or partly cloudy).

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

**Ans:** drop_first = True is used to drop the first index of every category. It is important during dummy variable creation because that was obvious result. So, it is better to drop it to manage space and time of the model.

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

**Ans:** temp and atemp seems to have highest correlation with the target variable at the pair_plot among the numerical variables.

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

**Ans:** By refer the map in a repetitive manner while validating linear model as to validate different correlated value along with p-values and VIFs, the assumptions of linear regression can be validated after building model on the training set.

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

**Ans:** Based on the final model, 'yr', 'mnth_9' and 'season_4' are the top 3 features contributing significantly towards explaining the demand of the shared bikes.

# General Subjective Questions

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

**Ans:** Linear regression is one of the very basic forms of machine learning where we train a model to predict the behavior of your data based on some variables. In the case of linear regression as you can see the name suggests linear that means the two variables which are on the x-axis and y-axis should be linearly correlated.

An example is let's say you are running a sales promotion and expecting a certain number of count of customers to be increased now what you can do is you can look the previous promotions and plot if over on the chart when you run it and then try to see whether there is an increment into the number of customers whenever you rate the promotions and with the help of the previous historical data you try to figure it out or you try to estimate what will be the count or what will be the estimated count for my current promotion this will give you an idea to do the planning in a much better way about how many numbers of stalls maybe you need or how many increase number of employees you need to serve the customer. Here the idea is to estimate the future value based on the historical data by learning the behavior or patterns from the historical data.

In some cases, the value will be linearly upward that means whenever X is increasing Y is also increasing or vice versa that means they have a correlation or there will be a linear downward relationship.

One example for that could be that the police department is running a campaign to reduce the number of robberies, in this case, the graph will be linearly downward.
Linear regression is used to predict a quantitative response Y from the predictor variable X.

Mathematically, we can write a linear regression equation as:   **y = a + bx**
Where a and b given by the formulas:

$$b\,(slope) = \frac{n \sum xy - \left(\sum x\right)\left(\sum y\right)}{n \sum x^2 - \left(\sum x\right)^2}$$

$$a\,(intercept) = \frac{n \sum y - b\left(\sum x\right)}{n}$$

Here, x and y are two variables on the regression line.
b = Slope of the line.
a = y-intercept of the line.
x = Independent variable from dataset
y = Dependent variable from dataset

Use Cases of Linear Regression:

1. Prediction of trends and Sales targets – To predict how industry is performing or how many sales targets industry may achieve in the future.
2. Price Prediction – Using regression to predict the change in price of stock or product.
3. Risk Management- Using regression to the analysis of Risk Management in the financial and insurance sector.

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

**Ans:** Anscombe's Quartet can be defined as a group of 4 data sets which are nearly identical in simple descriptive statistics, but there are some peculiarities in the dataset that fools the regression model if built. They have very different distributions and appear differently when plotted on scatter plots.

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

**Ans:** In Statistics, the Pearson's Correlation Coefficient is also referred to as **Pearson's r, the Pearson product-moment correlation coefficient (PPMCC), or bivariate correlation**. It is a statistic that measures the linear correlation between two variables. Like all correlations, it also has a numerical value that lies between -1.0 and +1.0.

Whenever we discuss correlation in statistics, it is generally Pearson's correlation coefficient. However, it cannot capture nonlinear relationships between two variables and cannot differentiate between dependent and independent variables.

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

**Ans:** Scaling is the procedure of measuring and assigning the objects to the numbers according to the specified rules. In other words, the process of locating the measured objects on the continuum, a continuous sequence of numbers to which the objects are assigned is called as scaling.

Feature scaling is a method used to normalize the range of independent variables or features of data. In data processing, it is also known as data normalization and is generally performed during the data pre-processing step.

The two most discussed scaling methods are Normalization and Standardization. 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).

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

**Ans:** If there is a perfect correlation, then the value of VIF is infinite.

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

**Ans:** A Q–Q (quantile-quantile) plot is a probability plot, which is a graphical method for comparing two probability distributions by plotting their quantiles against each other.

A Q–Q plot is used to compare the shapes of distributions, providing a graphical view of how properties such as location, scale, and skewness are similar or different in the two distributions.

This helps in a scenario of linear regression when we have training and test data set received separately and then we can confirm using Q-Q plot that both the data sets are from populations with same distributions.