# 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 –**
   a) Spring Season attracts less customer. While in fall season bike sharing has highest demand.
   b) Mean demand is almost same in the months of June, July, Aug, Sept but Sept has more demand when compared.
   c) Demand has picked-up in the year 2019
   d) Demand is high when weather is clear which is quite obvious.

2. Why is it important to use **drop_first=True** during dummy variable creation?
   **Ans –**
   "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?
   **Ans –**
   'temp' has highest correlation with the target variable.

4. How did you validate the assumptions of Linear Regression after building the model on the training set?
   **Ans –**
   a) Linear Relationship
   b) Homoscedasticity
   c) Absence of Multicollinearity
   d) Independence of residuals (absence of auto-correlation)
   e) Normality of Errors

5. Based on the final model, which are the top 3 features contributing significantly towards explaining the demand of the shared bikes?
   **Ans –**
   temp (0.46 times), season_winter (0.29 times), mnth_sep (0.27 times)

# General Subjective Questions

1. Explain the linear regression algorithm in detail.

   **Ans –**

   Linear Regression Algorithm is a machine learning algorithm based on supervised learning. Linear regression is one of the very basic forms of machine learning where we train a model to predict the behaviour 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 behaviour 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:
a)  Prediction of trends and Sales targets – To predict how industry is performing or how many sales targets industry may achieve in the future.
b)  Price Prediction – Using regression to predict the change in price of stock or product.
c)  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 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.
    The quartet is often used to illustrate the importance of looking at a set of data graphically before starting to analyse according to a particular type of relationship, and the inadequacy of basic statistic properties for describing realistic datasets.

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 l
    Pearson's correlation coefficient is the covariance of the two variables divided by the product of their standard deviations. The form of the definition involves a "product moment", that is, the mean (the first moment about the origin) of the product of the mean-adjusted random variables; hence the modifier product-moment in the name.
    Pearson's Correlation Coefficient is named after Karl Pearson. He formulated the correlation coefficient from a related idea by Francis Galton in the 1880s.

4.  What is scaling? Why is scaling performed? What is the difference between normalized scaling and standardized scaling?
    **Ans –**
    a)  Scaling 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.
    b)  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.

c) Normalized scaling brings all of the data in the range of 0 and 1. sklearn.preprocessing.MinMaxScaler helps to implement normalization in python.

Standardization replaces the values by their Z scores. It brings all of the data into a standard normal distribution which has mean (μ) zero and standard deviation one (σ). sklearn.preprocessing.scale helps to implement standardization in python.

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

   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.
   **Ans –**

   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.

   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.