# 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)**

   Bike demand in the fall is the highest.
   Bike demand takes a dip in spring.
   Bike demand in year 2019 is higher as compared to 2018.
   Bike demand is high in the months from May to October.
   Bike demand is high if weather is clear or with mist cloudy while it is low when there is light rain or light snow.
   The demand of bike is almost similar throughout the weekdays.
   Bike demand doesn't change whether day is working day or not.

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

   It is important in order to achieve k-1 dummy variables as it can be used to delete extra column while creating dummy variables.

   For Example: We have three variables: Furnished, Semi-furnished and un-furnished. We can only take 2 variables as furnished will be 1-0, semi-furnished will be 0-1, so we don't need unfurnished as we know 0-0 will indicate un-furnished. So we can remove it
   It is also used to reduce the collinearity between dummy variables

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

   atemp and temp both have same correlation with target variable of 0.63 which is the highest among all numerical variables

4. **How did you validate the assumptions of Linear Regression after building the model on the training set?** **(3 marks)**

   Residuals distribution should follow normal distribution and centred around 0(mean = 0). We validate this assumption about residuals by plotting a distplot of residuals and see if residuals are following normal distribution or not. The above diagram shows that the residuals are distributed about mean = 0.

5.  **Based on the final model, which are the top 3 features contributing significantly towards explaining the demand of the shared bikes?** **(2 marks)**

    The top 3 features (rounded to 3 decimal places) are 1. temp - coefficient: 0.492 2. yr - coefficient: 0.233 3. weathersit_Light Snow & Rain – coefficient: -0.28

## General Subjective Questions

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

    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.

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

    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.
    **Simple understanding:**

    Once Francis John "Frank" Anscombe who was a statistician of great repute found 4 sets of 11 data-points in his dream and requested the council as his last wish to plot those points. Those 4 sets of 11 data-points are given below.

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

    Pearson correlation coefficient or Pearson's correlation coefficient or Pearson's r is defined in statistics as the measurement of the strength of the relationship between two variables and their association with each other.

In simple words, Pearson's correlation coefficient calculates the effect of change in one variable when the other variable changes.

For example: Up till a certain age, (in most cases) a child's height will keep increasing as his/her age increases. Of course, his/her growth depends upon various factors like genes, location, diet, lifestyle, etc.

This approach is based on covariance and thus is the best method to measure the relationship between two variables.

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

> 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.

- It brings all of the data in the range of 0 and 1 sklearn.preprocessing.MinMaxScaler helps to implement normalization in python. MinMaxScaling = x-min(x)/max(x)-min(x)

- 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). Standadisation: x=x-mean(x)/sd(x)

- sklearn.preprocessing.scale helps to implement standardization in python.One disadvantage of normalization over standardization is that it losessome information in the data, especially about outliers.

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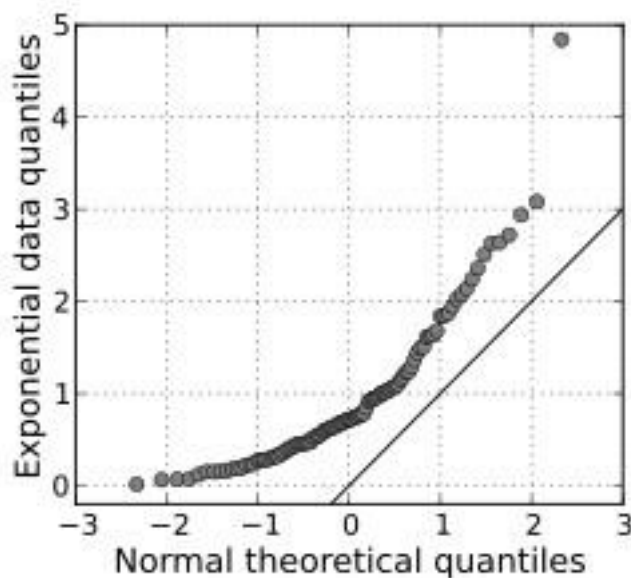**(3 marks)**

> 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.**

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.
A Q-Q plot showing the 45 degree 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.

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.