# 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?**

A:  From the given data, it gives me the understanding that the booking of bikes increased from 2019 than in 2018. And also from the date ranges in month, there was a increased in booking from the number ranges of ( 5, 6, 7 ,8 ,9)

--------------------------------------------------------------------------------------------------------------------------

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

A:  Because we can eliminate dummy variables through this function.

--------------------------------------------------------------------------------------------------------------------------

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

A:  'Atemp" has the highest correlation with target variable.

--------------------------------------------------------------------------------------------------------------------------

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

A: By using Linear regression, checked the linear correlation between variables such as cnt, months, seasons.

--------------------------------------------------------------------------------------------------------------------------

5. **Based on the final model, which are the top 3 features contributing significantly towards explaining the demand of the shared bikes?**
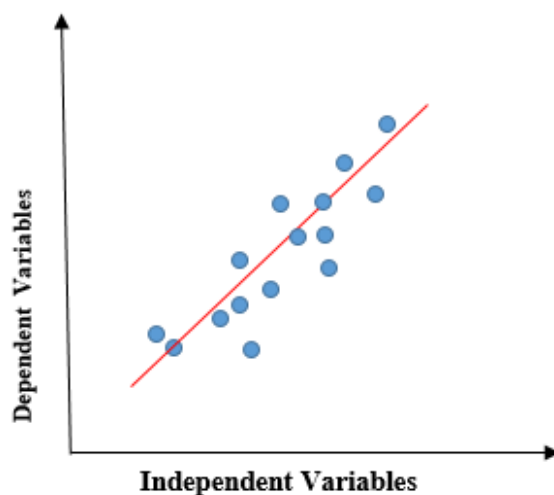
A: As mentioned, the top 3 contributing factors are that the demand was high between months ranging from 5-9 including weekdays.

--------------------------------------------------------------------------------------------------------------------------

## General questions:

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

A: Linear regression is a quiet and simple statistical regression method used for predictive analysis and shows the relationship between the continuous variables. Linear regression shows the linear relationship between the independent variable (X-axis) and the dependent

variable (Y-axis), consequently called linear regression. *If there is a single input variable (x), such linear regression is called simple linear regression. And if there is more than one input variable, such linear regression is called multiple linear regression.* The linear regression model gives a sloped straight line describing the relationship within the variables.



**Independent Variables**

The above graph presents the linear relationship between the dependent variable and independent variables. When the value of x (independent variable) increases, the value of y (dependent variable) is likewise increasing. The red line is referred to as the best fit straight line. Based on the given data points, we try to plot a line that models the points the best.

*To calculate best-fit line linear regression uses a traditional slope-intercept form.*

$$y = mx + b \implies y = a_0 + a_1 x$$

y= Dependent Variable.

x= Independent Variable.

a0= intercept of the line.

a1 = Linear regression coefficient.

-----------------------------------------------------------------------------------------------------------------

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

A: Anscombe's Quartet was developed by statistician Francis Anscombe. It comprises four datasets, each containing eleven (x, y) pairs. The essential thing to note about these datasets is that they share the same descriptive statistics. But things change completely, and I must emphasize COMPLETELY, when they are graphed. Each graph tells a different story irrespective of their similar summary statistics.

The summary statistics show that the means and the variances were identical for x and y across the groups:
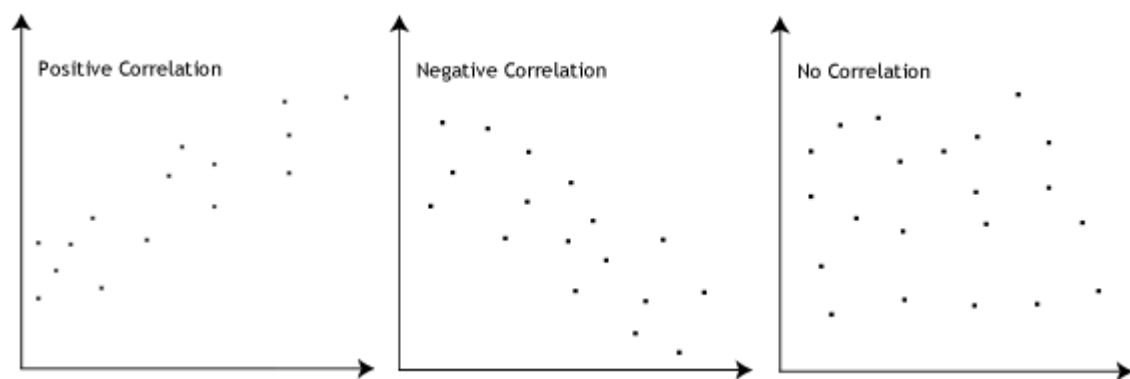
 • Mean of x is 9 and mean of y is 7.50 for each dataset.

 • Similarly, the variance of x is 11 and variance of y is 4.13 for each dataset

• The correlation coefficient (how strong a relationship is between two variables) between x and y is 0.816 for each dataset. Dataset I appears to have clean and well-fitting linear models.

• Dataset II is not distributed normally.

 • In Dataset III the distribution is linear, but the calculated regression is thrown off by an outlier. • Dataset IV shows that one outlier is enough to produce a high correlation coefficient. This quartet emphasizes the importance of visualization in Data Analysis. Looking at the data reveals a lot of the structure and a clear picture of the dataset.

-----------------------------------------------------------------------------------------------------------------

## 3. What is Pearson's R?

A: 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. The Pearson correlation coefficient, r, can take a range of values from +1 to -1. A value of 0 indicates that there is no association between the two variables. A value greater than 0 indicates a positive association; that is, as the value of one variable increases, so does the value of the other variable. A value less than 0 indicates a negative association; that is, as the value of one variable increases, the value of the other variable decreases. This is shown in the diagram below:

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



## Pearson r Formula

$$r = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2 \sum (y_i - \bar{y})^2}}$$

Here,

- =correlation coefficient

- =values of the x-variable in a sample

- =mean of the values of the x-variable

- =values of the y-variable in a sample

- =mean of the values of the y-variable

------------------------------------------------------------------------------

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

A: *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.*

$$\text{MinMax Scaling: } x = \frac{x - min(x)}{max(x) - min(x)}$$

*Standardization Scaling:*

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

$$\text{Standardisation: } x = \frac{x - mean(x)}{sd(x)}$$

- ***sklearn.preprocessing.scale** helps to implement standardization in python.*

- *One disadvantage of normalization over standardization is that it **loses** some information in the data, especially about **outliers**.*
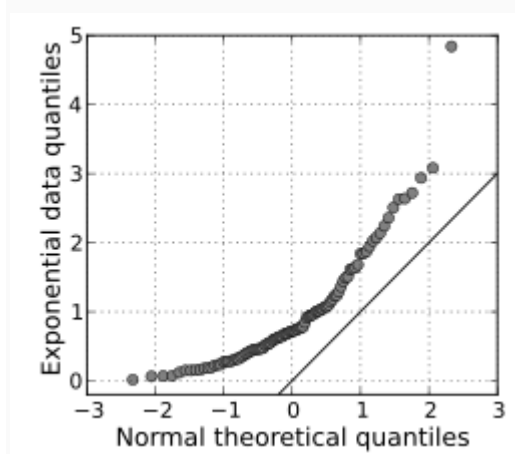
***Example:***

*Below shows example of Standardized and Normalized scaling on original values.*

| Original | Standardized | Normalized |
|---------|-------------|-----------|
| 59 | 1.85 | 1.000 |
| 11 | -1.06 | 0.000 |
| 55 | 1.61 | 0.917 |
| 30 | 0.09 | 0.396 |
| 14 | -0.88 | 0.063 |
| 18 | -0.64 | 0.146 |
| 32 | 0.21 | 0.438 |
| 27 | -0.09 | 0.333 |
| 21 | -0.46 | 0.208 |
| 18 | -0.64 | 0.146 |



---------------------------------------------------------------------------------

## 5.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.

---------------------------------------------------------------------------------------------------------------

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

A: 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).