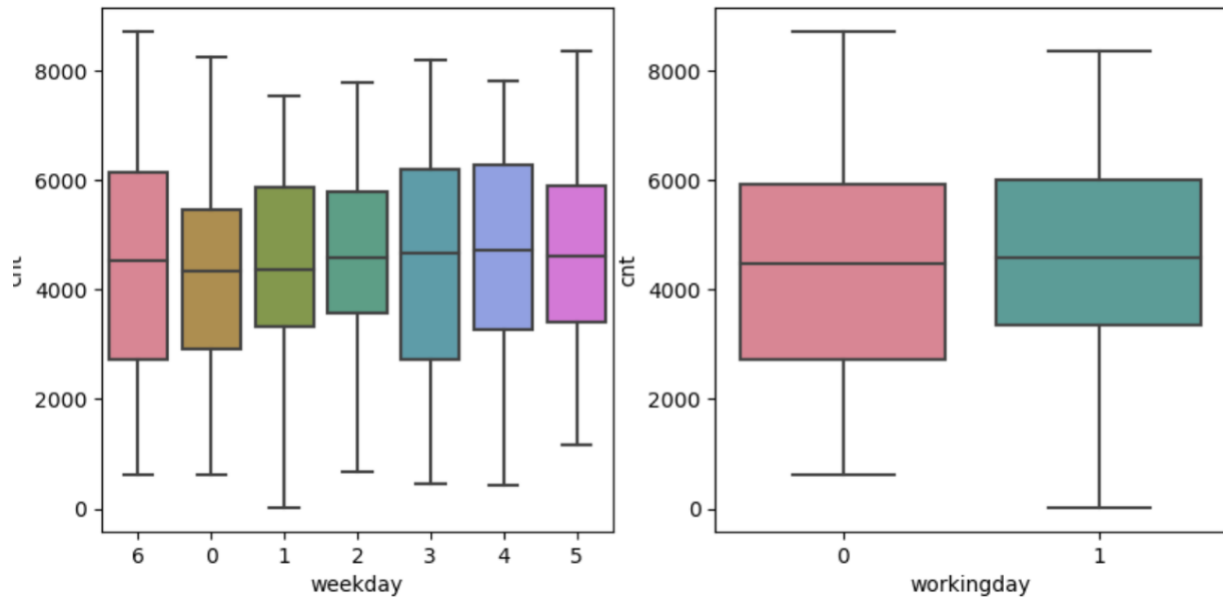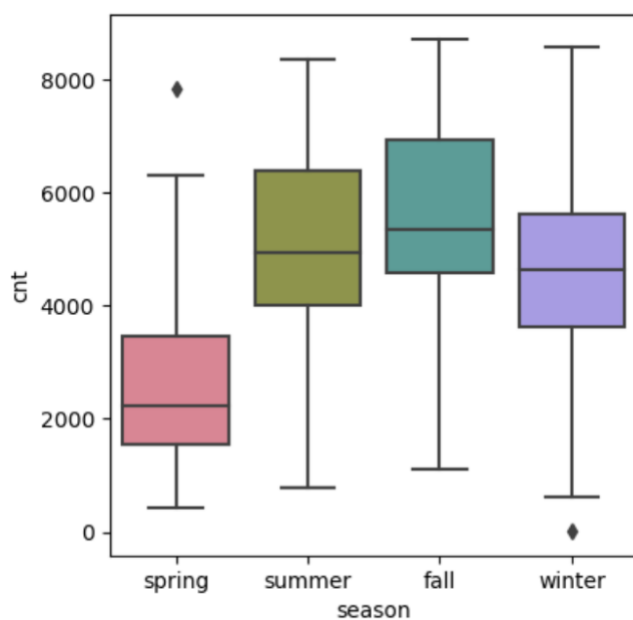# Assignment-based Subjective Questions

**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] Categorical variable are season, weekday and weathesit.



For the categorical variables weekday, working day there is consitent bookings and not must dependent, but when we have season, there is a considerable dependency on the bookings

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

[Answer] 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? (1 mark)**

[Answer] Temperature has the considerable corelation with target variable

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

[Answer] According to the made assumption there is linear relationship between the features and target. Linear regression captures only linear relationship. This can be validated by plotting a scatter plot between the features and the target.
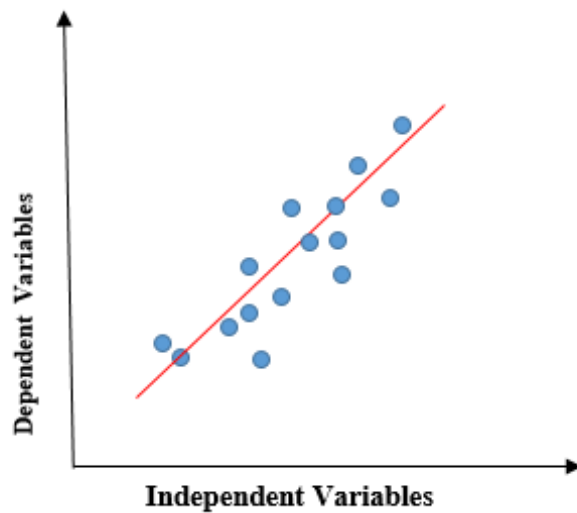
**5. Based on the final model, which are the top 3 features contributing significantly towards explaining the demand of the shared bikes?**
[Answer] Based on the assessment made using linear regression, the top features contributing are temparature, wheathersit.

# General Subjective Questions

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

**[Answer]** 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.

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.

$$y = \beta 0 + \beta 1x$$

Cost function is used to optimizes the regression coffiecients and measures how linear regression model is performing.
In Linear regression, Mean Squared Error cost function is used , which the avarage squared of difference occurred between actual and predicted.

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

[Answer] Anscombe's Quartet can be defined as a group of four 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.

It was constructed in 1973 by statistician Francis Anscombe to illustrate

the importance of plotting the graphs before analyzing and model building, and the effect of

other observations on statistical properties.There are these four data set plots which have

nearly same statistical observations, which provides same statistical information that

involves variance, and mean of all x,y points in all four datasets.

This tells us about the importance of visualising the data before applying various algorithms

out there to build models out of them which suggests that the data features must be plotted

in order to see the distribution of the samples that can help you identify the various

anomalies present in the data like outliers, diversity of the data, linear separability of the

data, etc. Also, the Linear Regression can be only be considered a fit for the data with linear

relationships and is incapable of handling any other kind of datasets.

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

[Answer] The **Pearson correlation coefficient** (*r*) is the most common way of measuring a linear correlation. It is a number between –1 and 1 that measures the strength and direction of the relationship between two variables.

If the cofficient lies between 0 and 1, positively corerelated
0 – No Corelation and -1 and 0 then Negatively correlated.

### 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.

Normalized/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.

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

Standardized 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 ($\sigma$).

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

- 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.

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

[Answer] If VIF = infinity which mean perfect correlation. 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.

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.