# Assignment-based Subjective Questions

Ans.1) From the given Dataset, we have- season, yr, mnth, holiday, weekday, workingday, weathersit.

cnt= 0.246 X yr - 0.083 X holiday -  0.198 X Spring - 0.321 X Light rain_Light snow_Thunderstorm - 0.090 X Mist_Cloudy +0.063 X 3 + 0.123  X 5 +0.148 X 6 +0.153 X 8 + 0.193 X 9 -0.049 X Sunday + 0.126 X 7 + 0.116 X 10

From the analysis we can see that these categorical variables affect the target variable which is cnt in the following ways:

- yr: has a positive coefficient of 0.2469 as per the summary and is significant because its p-value is almost 0. And demand is also directly proportional to the yr.
- mnth: From the final summary it is clear that demand for the cycles are higher in the month of 3,5,6,7,8,9 and 10
- weekday: Demand is low on Sundays.
- Holiday: Demand decreases on holidays.
- Season: Demand decreases when it is spring specifically.
- Weathersit: Demand decreases when it is Light_rain,Light_snow,Thunderstorm,Mist-Cloudy


Ans.2) "**drop_first=True**" is mainly used for dropping the first column when we create a dummy variable. It is done because whenever we create dummy variables for n categorical variables then they are represented by n-1 columns. It depends on the number of levels of the categorical variables, then they are represented by one less number of columns.

Ans.3) temp has the highest correlation with our target variable which is cnt. It is 63%

Ans.4) After building the model for Linear Regression we validated the assumptions by:

- Checking predicted value against the test data and implying that it shows a linear relationship which it should.
- Checking for the normal distribution of the errors by plotting a histogram for the error terms (y-ypred value).
- Calculating the R square for the test data and checking whether the model is still efficient or not.
- Checking the assumption of homoscedasticity. Whether residuals are showing a pattern or not by plotting a scatter graph.

Ans.5) Based on the model the top 3 features which contribute to the sales of the bikes are:

- Sales of the bikes increases in the month of 3,5,6,7,8,10
- Sales decreases if it is a holiday, so holidays contribute inversely.
- Sales also depend on seasons, the model shows that when there is a thunderstorm, rains, snow or misty; the sales drop.

**Ans.1) What is Linear Regression?**

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.

Example for that can be 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:

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.

Steps involving Linear Regression are:

- **Reading, understanding and visualizing the data**: These steps basically involve reading the data and looking for the null values, numerical and categorical variables. Visualizing the data using EDA techniques by plotting a scatterplot, pair plot or a heatmap.
- **Data Preparation**: In this step we convert the categorical variables into dummy variables in order to convert them to integer. Separate the categorical variables into a data frame and then merging them again once we have created the dummy variables.
- **Splitting the Data into Training and Test sets**: We use sklearn.model_selection module and import train_test_split() function in order to split the dataset into training and test set. We also do Rescaling in this step, it is an important step; because it is very important for the variables to have comparable scales. Otherwise, some of the coefficients obtained by fitting the regression model might be very large or very small compared to other coefficients. Scaling is done by using two methods: 1) Min-Max Scaling 2) Standardisation.
- **Divide the dependent and the predictor variables in X and y sets for the model building**.
- **Building a linear model**:  Fit a regression line through the training data using statsmodel api.
- **Feature Selection**: We have to select the variables which will increase the adjusted R-squared of the model in order to be an efficient one. We do so by comparing the p-value and the VIF(virtual inflation factor). For doing so we have mainly three approaches: 1) Manual approach 2) Automated approach 3) Mixed approach. For the sake of model building, we use the third one.
- **Residual Analysis on the train data**: We need to check whether the error terms are normally distributed or not and we do so by plotting the histogram.
- **Making Predictions using the final model**: Now we work with the test set and instead of using the fit() function we only use the transform() with the test set, because we don't want the model to learn about the test set in advance.
- **Model Evaluation**: Plot the final graph and determine the equation in order to find the predictors which determine the value of the dependent variable.

Ans.2) Anscombe's Quartet

According to the definition given in Wikipedia, **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:
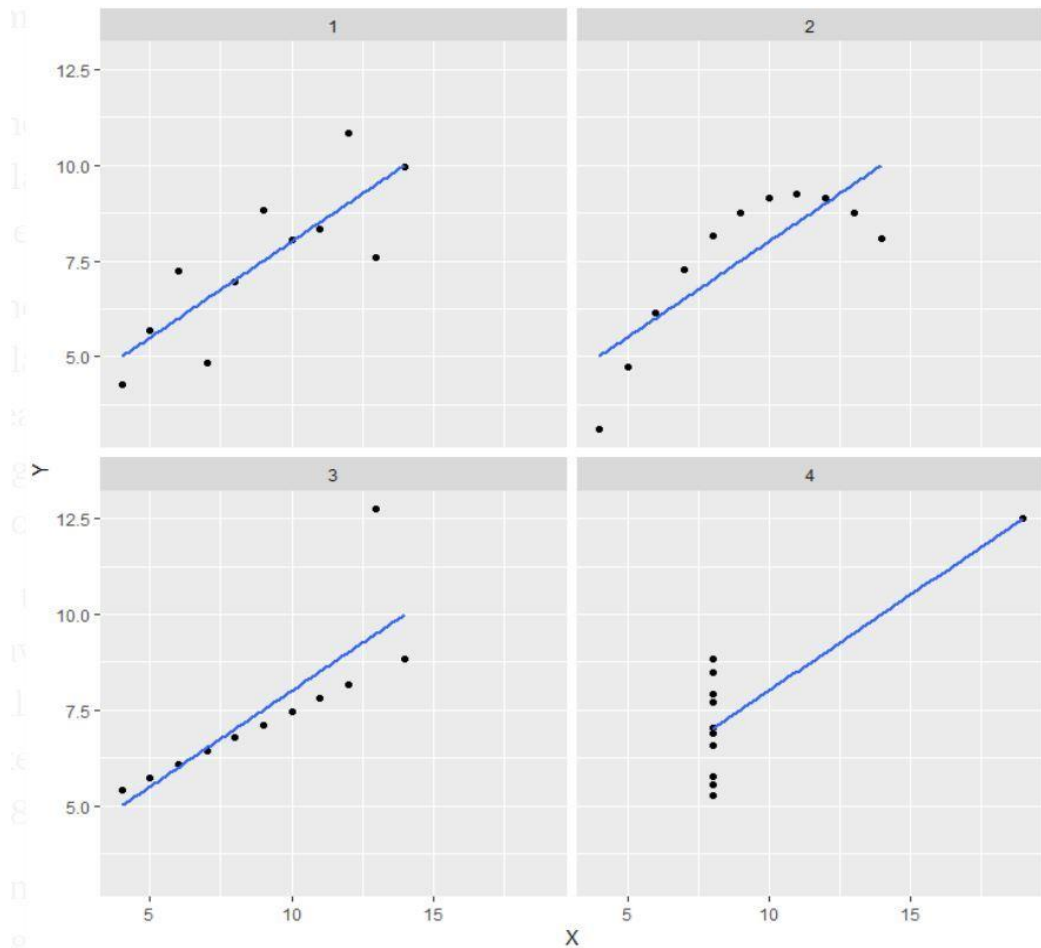
```
+-------+--------+-------+--------+-------+--------+-------+-------+
|      I         |      II        |      III        |      IV        |
+-------+--------+-------+--------+-------+--------+-------+-------+
| x     | y      | x     | y      | x     | y      | x     | y     |
----+--------+-------+-------+-------+-------+-------+-------+------+
| 10.0  | 8.04   | 10.0  | 9.14   | 10.0  | 7.46   | 8.0   | 6.58  |
| 8.0   | 6.95   | 8.0   | 8.14   | 8.0   | 6.77   | 8.0   | 5.76  |
| 13.0  | 7.58   | 13.0  | 8.74   | 13.0  | 12.74  | 8.0   | 7.71  |
| 9.0   | 8.81   | 9.0   | 8.77   | 9.0   | 7.11   | 8.0   | 8.84  |
| 11.0  | 8.33   | 11.0  | 9.26   | 11.0  | 7.81   | 8.0   | 8.47  |
| 14.0  | 9.96   | 14.0  | 8.10   | 14.0  | 8.84   | 8.0   | 7.04  |
| 6.0   | 7.24   | 6.0   | 6.13   | 6.0   | 6.08   | 8.0   | 5.25  |
| 4.0   | 4.26   | 4.0   | 3.10   | 4.0   | 5.39   | 19.0  | 12.50 |
| 12.0  | 10.84  | 12.0  | 9.13   | 12.0  | 8.15   | 8.0   | 5.56  |
| 7.0   | 4.82   | 7.0   | 7.26   | 7.0   | 6.42   | 8.0   | 7.91  |
| 5.0   | 5.68   | 5.0   | 4.74   | 5.0   | 5.73   | 8.0   | 6.89  |
+-------+--------+-------+--------+-------+--------+-------+-------+
```

After that, the council analyzed them using only descriptive statistics and found the mean, standard deviation, and correlation between x and y.

Summary for the above analysis:

```
                    Summary
+-----+---------+-------+---------+-------+----------+
| Set | mean(X) | sd(X) | mean(Y) | sd(Y) | cor(X,Y) |
+-----+---------+-------+---------+-------+----------+
|  1  |       9 | 3.32  |     7.5 | 2.03  |   0.816  |
|  2  |       9 | 3.32  |     7.5 | 2.03  |   0.816  |
|  3  |       9 | 3.32  |     7.5 | 2.03  |   0.816  |
|  4  |       9 | 3.32  |     7.5 | 2.03  |   0.817  |
+-----+---------+-------+---------+-------+----------+
```

Regression line plots:



**Explanation of this output:**

- In the first one(top left) if you look at the scatter plot you will see that there seems to be a linear relationship between x and y.

- In the second one(top right) if you look at this figure you can conclude that there is a non-linear relationship between x and y.

- In the third one(bottom left) you can say when there is a perfect linear relationship for all the data points except one which seems to be an outlier which is indicated be far away from that line.

- Finally, the fourth one(bottom right) shows an example when one high-leverage point is enough to produce a high correlation coefficient.

**Application:**
The quartet is still often used to illustrate the importance of looking at a set of data graphically before starting to analyze according to a particular type of relationship, and the inadequacy of basic statistic properties for describing realistic datasets.

Ans.3) **Correlation coefficients** are used to measure how strong a relationship is between two variables. There are several types of correlation coefficient, but the most popular is Pearson's. **Pearson's correlation** (also called Pearson's *R*) is a **correlation coefficient** commonly used in linear regression.

here are several types of correlation coefficient formulas.

One of the most commonly used formulas is Pearson's correlation coefficient formula. If you're taking a basic stats class, this is the one you'll probably use:

$$r = \frac{n(\sum xy) - (\sum x)(\sum y)}{\sqrt{[n\sum x^2 - (\sum x)^2][n\sum y^2 - (\sum y)^2]}}$$

Correlation between sets of data is a measure of how well they are related. The most common measure of correlation in stats is the Pearson Correlation. The full name is the **Pearson Product Moment Correlation (PPMC)**. It shows the linear relationship between two sets of data. In simple terms, it answers the question, *Can I draw a line graph to represent the data?* Two letters are used to represent the Pearson correlation: Greek letter rho ($\rho$) for a population and the letter "r" for a sample.

Potential problems with Pearson correlation.

The PPMC is not able to tell the difference between dependent variables and independent variables. For example, if you are trying to find the correlation between a high calorie diet and diabetes, you might find a high correlation of .8. However, you could also get the same result with the variables switched around. In other words, you could say that diabetes causes a high calorie diet. That obviously makes no sense. Therefore, as a researcher you have to be aware of the data you are plugging in. In addition, the PPMC will not give you any information about the slope of the line; it only tells you whether there is a relationship.

**Real Life Example**

Pearson correlation is used in thousands of real life situations. For example, scientists in China wanted to know if there was a relationship between how weedy rice populations are different genetically. The goal was to find out the evolutionary potential of the rice. Pearson's correlation between the two

groups was analyzed. It showed a positive Pearson Product Moment correlation of between 0.783 and 0.895 for weedy rice populations. This figure is quite high, which suggested a fairly strong relationship.

Ans.4) Scaling is done to transform all the variables into a single comparable scale. If we don't have comparable scales, then the coefficients obtained by fitting the regression model might be very large or very small as compared to other coefficients. This might become very annoying at the time of model evaluation. So it is advised to use standardization or normalization so that the units of the coefficients obtained are all on the same scale.

**What is Normalization?**

**Normalization is a scaling technique in which values are shifted and rescaled so that they end up ranging between 0 and 1. It is also known as Min-Max scaling.**

Here's the formula for normalization:

$$X' = \frac{X - X_{min}}{X_{max} - X_{min}}$$

Here, Xmax and Xmin are the maximum and the minimum values of the feature respectively.

- When the value of X is the minimum value in the column, the numerator will be 0, and hence X' is 0

- On the other hand, when the value of X is the maximum value in the column, the numerator is equal to the denominator and thus the value of X' is 1

- If the value of X is between the minimum and the maximum value, then the value of X' is between 0 and 1

**What is Standardization?**

**Standardization is another scaling technique where the values are centered around the mean with a unit standard deviation. This means that the mean of the attribute becomes zero and the resultant distribution has a unit standard deviation.**

Here's the formula for standardization:

$$X' = \frac{X - \mu}{\sigma}$$

$\mu$ is the mean of the feature values and $\sigma$ is the standard deviation of the feature values. Note that in this case, the values are not restricted to a particular range.

Normalization vs. standardization is an eternal question among machine learning newcomers. Let me elaborate on the answer in this section.

- Normalization is good to use when you know that the distribution of your data does not follow a Gaussian distribution. This can be useful in algorithms that do not assume any distribution of the data like K-Nearest Neighbors and Neural Networks.

- Standardization, on the other hand, can be helpful in cases where the data follows a Gaussian distribution. However, this does not have to be necessarily true. Also, unlike normalization, standardization does not have a bounding range. So, even if you have outliers in your data, they will not be affected by standardization.

Ans.5) What is VIF?

VIF is an index that provides a measure of how much the variance of an estimated regression coefficient increases due to collinearity. In order to determine VIF, we fit a regression model between the independent variables. For example, we would fit the following models to estimate the coefficient of determination R1 and use this value to estimate the VIF:

$X\_1 = C + \alpha\_2 X\_2 + \alpha\_3 X\_3 + \cdots$

$〚VIF〛\_1 = 1/(1 - R\_1^2)$

Next, we fit the model between X2 and the other independent variables to estimate the coefficient of determination R2:

$X\_2 = C + \alpha\_1 X\_1 + \alpha\_3 X\_3 + \cdots$

$〚VIF〛\_2 = 1/(1 - R\_2^2)$

If all the independent variables are orthogonal to each other, then VIF = 1.0. If there is perfect correlation, then VIF = infinity. A large value of VIF indicates that there is a correlation between the variables. If the VIF is 4, this means that the variance of the model coefficient is inflated by a factor of 4 due to the presence of multicollinearity. This would mean that that standard error of this coefficient is inflated by a factor of 2 (square root of variance is the standard deviation). The standard error of the coefficient determines the confidence interval of the model coefficients. If the standard error is large, then the confidence intervals may be large, and the model coefficient may come out to be non-significant due to the presence of multicollinearity. A general rule of thumb is that if VIF > 10 then there is multicollinearity. Note that this is a rough rule of thumb, in some cases we might choose to live with

high VIF values if it does not affect our model results such as when we are fitting a quadratic or cubic model or depending on the sample size a large value of VIF may not necessarily indicate a poor model.

If there is a perfect correlation between two variables then the VIF values is infinity.

Ans.6) Quantile-Quantile (Q-Q) plot, is a graphical tool to help us assess if a set of data plausibly came from some theoretical distribution such as a Normal, exponential or Uniform distribution. Also, it helps to determine if two data sets come from populations with a common distribution.

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.

**Few advantages:**

a) It can be used with sample sizes also

b) Many distributional aspects like shifts in location, shifts in scale, changes in symmetry, and the presence of outliers can all be detected from this plot.

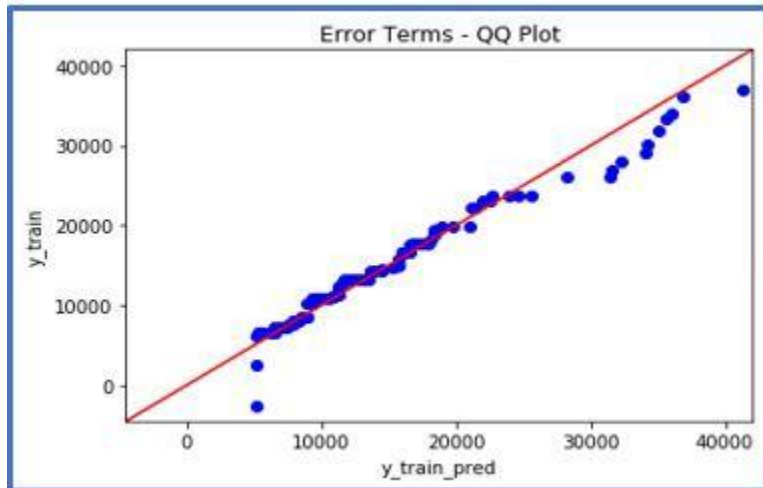It is used to check following scenarios:

If two data sets —

i. come from populations with a common distribution

ii. have common location and scale

iii. have similar distributional shapes
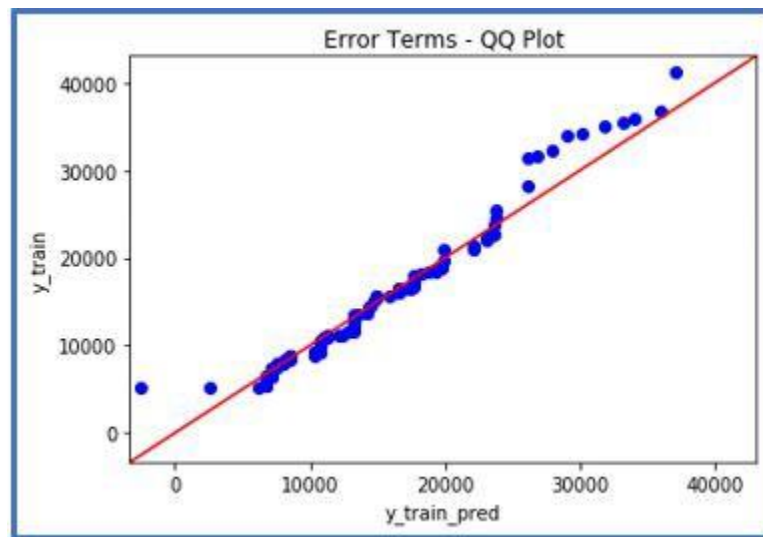
iv. have similar tail behavior

**Interpretation:**

A q-q plot is a plot of the quantiles of the first data set against the quantiles of the second data set.

Below are the possible interpretations for two data sets.

a) **Similar distribution**: If all point of quantiles lies on or close to straight line at an angle of 45 degree from x -axis

b) **Y-values < X-values:** If y-quantiles are lower than the x-quantiles.

**c) X-values < Y-values:** If x-quantiles are lower than the y-quantiles.



**d) Different distribution:** If all point of quantiles lies away from the straight line at an angle of 45 degree from x -axis