

## Assignment-based Subjective Questions

**Q1. From your analysis of the categorical variables from the dataset, what could you infer about their effect on the dependent variable?**

**Ans:** As per the analysis of data set, we can see the following observations for each categorical columns

**Season:**

- Maximum demand- Fall Season
- Moderate demand - Winter and summer season
- Lowest demand - Spring season

**Holiday:** Demand is higher in the holidays compared to weekdays

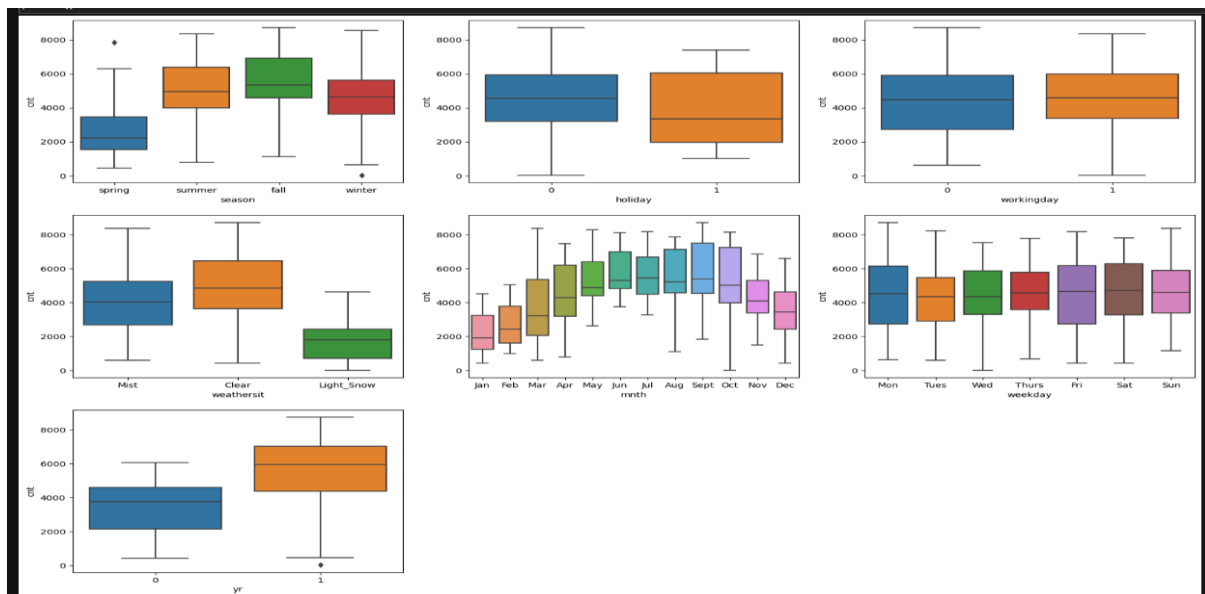
**Working day:** Demand is higher in the non-working day.

**Weathersit:** Demand is higher in the case of Clear weather, Few clouds, Partly cloudy and Partly cloudy lowest in the Light Snow, Light Rain + Thunderstorm + Scattered clouds and Light Rain + Scattered clouds. There's no record for Heavy Rain + Ice Pellets + Thunderstorm + Mist, Snow + Fog, so no conclusion can be made.

**Month:** There is gradually increase in demand from first quarter to third quarter and decrease from high to low in the last quarter of the year.

**Year:** There is in demand increase as the year passes.

Please get more details in following screen.



**Q2. 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.

**Q3. Looking at the pair-plot among the numerical variables, which one has the highest correlation with the target variable?**

**Ans:** Registered (95% Correlation)

**Q4 How did you validate the assumptions of Linear Regression after building the model on the training set?**

**Ans:**

- Verified the p value for each feature to know the significance.
- Checked the VIF value for each feature present in the data model.
- Calculated the R-square value on the test data set based on trained data model, which is approx. 78% (Good).
- By plotting the graph between y\_test\_pred and y\_test.
- Error terms is normalised in nature and the peak at 0.

**Q5. Based on the final model, which are the top 3 features contributing significantly towards explaining the demand of the shared bikes?**

**Ans:-** Temperature, weathesit and year

## General Subjective Questions

**Q1. Explain the linear regression algorithm in detail.**

**Ans:** 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.

Mathematically, we can write a linear regression equation as:

$$y = a + bx$$

**Where a and b given by the formulas:**

$$b(\text{slope}) = \frac{n \sum xy - (\sum x)(\sum y)}{n \sum x^2 - (\sum x)^2} \quad a(\text{intercept}) = \frac{n \sum y - b(\sum x)}{n}$$

**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

## Q2: Explain the Anscombe's quartet in detail.

**Ans:** 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 analysing 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.

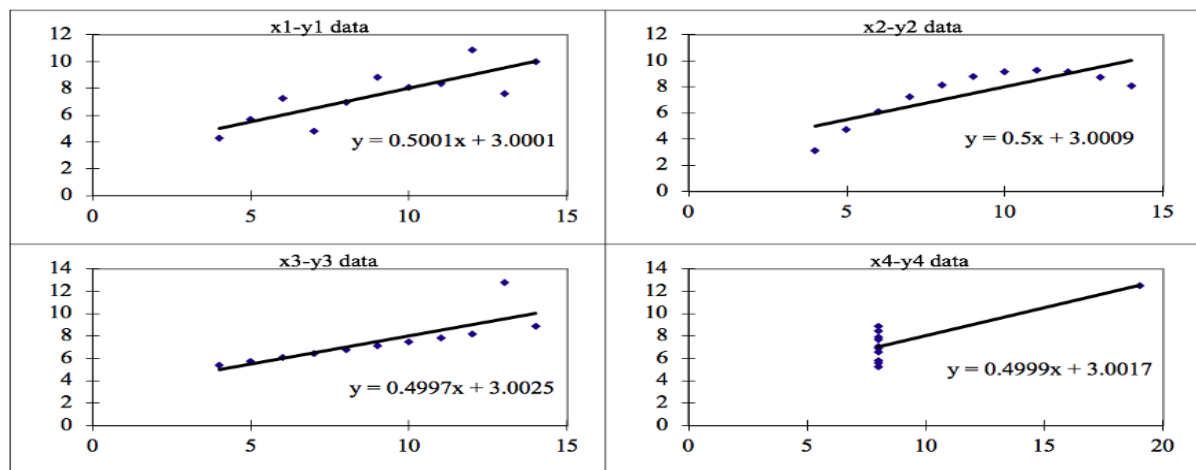
This tells us about the importance of visualising the data before applying various algorithms out there to build models out of them which suggests that the data features must be plotted in order to see the distribution of the samples that can help you identify the various anomalies present in the data 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. These four plots can be defined as follows:

Anscombe's Data											
Observation	x1	y1		x2	y2		x3	y3		x4	y4
1	10	8.04		10	9.14		10	7.46		8	6.58
2	8	6.95		8	8.14		8	6.77		8	5.76
3	13	7.58		13	8.74		13	12.74		8	7.71
4	9	8.81		9	8.77		9	7.11		8	8.84
5	11	8.33		11	9.26		11	7.81		8	8.47
6	14	9.96		14	8.1		14	8.84		8	7.04
7	6	7.24		6	6.13		6	6.08		8	5.25
8	4	4.26		4	3.1		4	5.39		19	12.5
9	12	10.84		12	9.13		12	8.15		8	5.56
10	7	4.82		7	7.26		7	6.42		8	7.91
11	5	5.68		5	4.74		5	5.73		8	6.89

The statistical information for all these four datasets are approximately similar and can be computed as follows:

Anscombe's Data											
Observation	x1	y1		x2	y2		x3	y3		x4	y4
1	10	8.04		10	9.14		10	7.46		8	6.58
2	8	6.95		8	8.14		8	6.77		8	5.76
3	13	7.58		13	8.74		13	12.74		8	7.71
4	9	8.81		9	8.77		9	7.11		8	8.84
5	11	8.33		11	9.26		11	7.81		8	8.47
6	14	9.96		14	8.1		14	8.84		8	7.04
7	6	7.24		6	6.13		6	6.08		8	5.25
8	4	4.26		4	3.1		4	5.39		19	12.5
9	12	10.84		12	9.13		12	8.15		8	5.56
10	7	4.82		7	7.26		7	6.42		8	7.91
11	5	5.68		5	4.74		5	5.73		8	6.89
				Summary Statistics							
N	11	11		11	11		11	11		11	11
mean	9.00	7.50		9.00	7.500909		9.00	7.50		9.00	7.50
SD	3.16	1.94		3.16	1.94		3.16	1.94		3.16	1.94
r	0.82			0.82			0.82			0.82	

When these models are plotted on a scatter plot, all datasets generate a different kind of plot that is not interpretable by any regression algorithm which is fooled by these peculiarities and can be seen as follows:



The four datasets can be described as:

**Dataset 1:** this fits the linear regression model well.

**Dataset 2:** this could not fit linear regression model on the data quite well as the data is non-linear.

**Dataset 3:** shows the outliers involved in the dataset which cannot be handled by linear regression model

**Dataset 4:** shows the outliers involved in the dataset which cannot be handled by linear regression model

**Q3: What is Pearson's R?**

**Ans:** 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.

**For example:**

**Positive linear relationship:** In most cases, universally, the income of a person increases as his/her age increases.

**Negative linear relationship:** If the vehicle increases its speed, the time taken to travel decreases, and vice versa.

From the example above, it is evident that the Pearson correlation coefficient,  $r$ , tries to find out two things – the strength and the direction of the relationship from the given sample sizes.

#### **Pearson correlation coefficient formula**

The correlation coefficient formula finds out the relation between the variables. It returns the values between -1 and 1. Use the below Pearson coefficient correlation calculator to measure the strength of two variables.

**Pearson correlation coefficient formula:**

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

Where:

$N$  = the number of pairs of scores

$\sum xy$  = the sum of the products of paired scores

$\sum x$  = the sum of  $x$  scores

$\sum y$  = the sum of  $y$  scores

$\sum x^2$  = the sum of squared  $x$  scores

$\sum y^2$  = the sum of squared  $y$  scores

**Q4. What is scaling? Why is scaling performed? What is the difference between normalized scaling and standardized scaling?**

**Ans:**

**What?**

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.

**Why?**

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 must 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 ( $\mu$ ) zero and standard deviation one ( $\sigma$ ).

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

*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.

**Q5. You might have observed that sometimes the value of VIF is infinite. Why does this happen?**

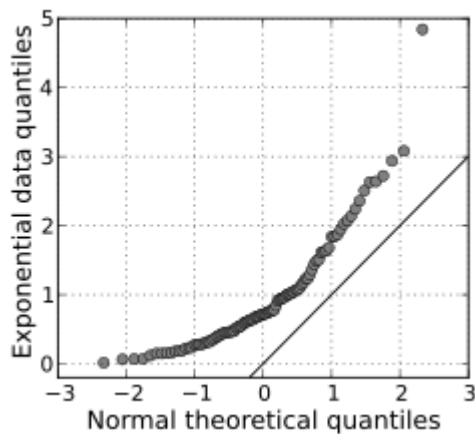
**Ans:** If there is perfect correlation, then VIF = infinity. This shows a perfect correlation between two independent variables. In the case of perfect correlation, we get  $R^2 = 1$ , which lead to  $1/(1-R^2)$  infinity. To solve this problem, we need to drop one of the variables from the dataset 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).

**Q6. 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. 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.