

## **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)**

- Season column is categorical column with each value corresponding to specific season.  
Ride Count Seems to be in maximum in Fall (Autumn) followed by Summer, Spring & Winter respectively.
- Year column has 2 values 0 and 1 and is a categorical column  
Ride Count has increased drastically in 2019 as compared to 2018
- Month column is categorical column with each value corresponding to specific month  
Ride Count seems to increase between May to October which are comparatively Fall(Autumn) & Summer Season in US
- Holiday/Weekday is a Categorical Variable.  
Ride Count is lesser on Holidays as compared to other days.i.e.,more bikes are rented during normal working days than on weekends or holidays.Weather Situation is a Categorical Variable with following values  
  
1: Clear, Few clouds, partly cloudy, partly cloudy  
2: Mist + Cloudy, Mist + Broken clouds, Mist + Few clouds, Mist  
3: Light Snow, Light Rain + Thunderstorm + Scattered clouds, Light Rain + Scattered clouds  
4: Heavy Rain + Ice Pallets + Thunderstorm + Mist, Snow + Fog  
  
Ride Count is more on Clear & Misty Days as compared to Light Snow / Rainfall.  
More bikes are rent during Clear, Few clouds, partly cloudy weather.

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

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.

Let's say we have 3 types of values in Categorical column and we want to create dummy variable for that column. If one variable is not furnished and semi\_furnished, then it is obvious unfurnished. So, we do not need 3<sup>rd</sup> variable to identify the unfurnished. Hence if we have categorical variable with n-levels, then we need to use n-1 columns to represent the dummy variables.

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

By looking at the pair plot temp variable has the highest (0.63) correlation with target variable 'cnt'.

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

The simple way to determine if this assumption is met or not is by creating a scatter plot x vs y. If the data points fall on a straight line in the graph, there is a linear relationship between the dependent and the independent variables, and the assumption holds.

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

Answer: The Top 3 features contributing significantly towards the demands of share bikes are:

- ☐ weathersit\_Light\_Snow(negative correlation).
- ☐ yr\_2019(Positive correlation).
- ☐ temp(Positive correlation).

### **General Subjective Questions**

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

Linear regression may be defined as the statistical model that analyzes the linear relationship between a dependent variable with given set of independent variables. Linear relationship between variables means that when the value of one or more independent variables will change (increase or decrease), the value of dependent variable will also change accordingly (increase or decrease).

Mathematically the relationship can be represented with the help of following equation –

$$Y=mX+b$$

Here, Y is the dependent variable we are trying to predict.

X is the independent variable we are using to make predictions.

m is the slop of the regression line which represents the effect X has on Y

b is a constant, known as the Y-intercept. If  $X = 0$ , Y would be equal to b.

Furthermore, the linear relationship can be positive or negative in nature as explained below –

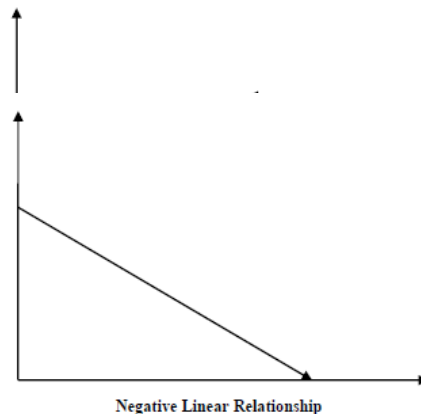
#### **Positive Linear Relationship**

A linear relationship will be called positive if both independent and dependent variable increases. It can be understood with the help of following graph –

#### **Negative Linear relationship**

A linear relationship will be called positive if independent increases and dependent variable increases. It can be understood with the help of following graph –

Assumptions



The following are some assumptions about dataset that is made by Linear Regression model –

**Multi-collinearity** – Linear regression model assumes that there is very little or no multi-collinearity in the data. Basically, multi-collinearity occurs when the independent variables or features have dependency in them.

**Auto-correlation** – Another assumption Linear regression model assumes is that there is very little or no auto-correlation in the data. Basically, auto-correlation occurs when there is dependency between residual errors.

**Relationship between variables** – Linear regression model assumes that the relationship between response and feature variables must be linear.

## 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.

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.

Code: Python program to find mean, standard deviation, and the correlation between x and y

```
# Import the required libraries
```

```
import pandas as pd
```

```
import statistics
```

```
from scipy.stats import pearsonr
```

```
# Import the csv file
```

```
df = pd.read_csv("anscombe.csv")
```

```
# Convert pandas dataframe into pandas series
```

```
list1 = df['x1']
```

```
list2 = df['y1']
```

```
# Calculating mean for x1
```

```
print('%.1f' % statistics.mean(list1))
```

```
# Calculating standard deviation for x1
```

```
print('%.2f' % statistics.stdev(list1))
```

```
# Calculating mean for y1
```

```
print('%.1f' % statistics.mean(list2))
```

```
# Calculating standard deviation for y1
```

```
print('%.2f' % statistics.stdev(list2))
```

```
# Calculating pearson correlation
```

```
corr, _ = pearsonr(list1, list2)
print('%.3f % corr)
```

# Similarly calculate for the other 3 samples

# This code is contributed by Amiya Rout

Output:

9.0

3.32

7.5

2.03

0.816

So let me show you the result in a tabular fashion for better understanding.

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	

Code: Python program to plot scatter plot

# Import the required libraries

from matplotlib import pyplot as plt

import pandas as pd

# Import the csv file

df = pd.read\_csv("anscombe.csv")

```
# Convert pandas dataframe into pandas series
```

```
list1 = df['x1']
```

```
list2 = df['y1']
```

```
# Function to plot scatter
```

```
plt.scatter(list1, list2)
```

```
# Function to show the plot
```

```
plt.show()
```

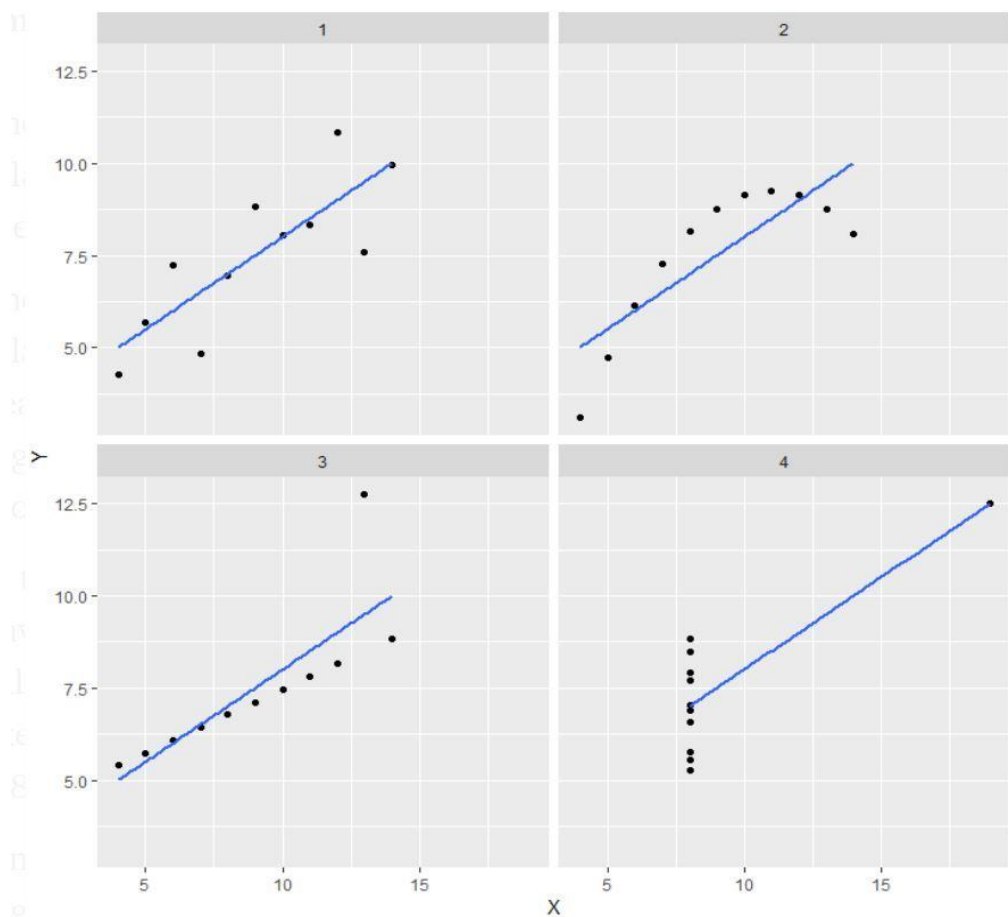
```
# Similarly plot scatter plot for other 3 data sets
```

For regression

line

refer

Output:



Note: It is mentioned in the definition that Anscombe's quartet comprises four datasets that have nearly identical simple statistical properties, yet appear very different when graphed.

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.

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

The Pearson correlation coefficient is used to measure the strength of a linear association between two variables, where the value  $r = 1$  means a perfect positive correlation and the value  $r = -1$  means a perfect negative correlation. So, for example, you could use this test to find out whether people's height and weight are correlated (the taller the people are, the heavier they're likely to be).

Requirements for Pearson's correlation coefficient are as follows:Scale of measurement should be interval or ratio

- ☐ Variables should be approximately normally distributed
- ☐ The association should be linear
- ☐ There should be no outliers in the data

Equation

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

What does this test do?

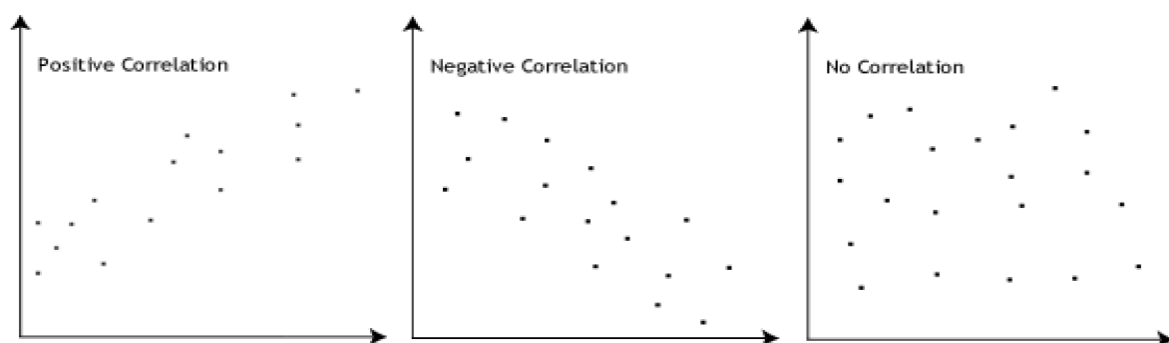
The Pearson product-moment correlation coefficient (or Pearson correlation coefficient, for short) is a measure of the strength of a linear association between two variables and is denoted by 'r'. Basically, a Pearson product-moment correlation attempts to draw a line of best fit through the data of two variables, and the Pearson correlation coefficient, r, indicates how far away all these data points are to this line of best fit (i.e., how well the data points fit this new model/line of best fit).

What values can the Pearson correlation coefficient take?

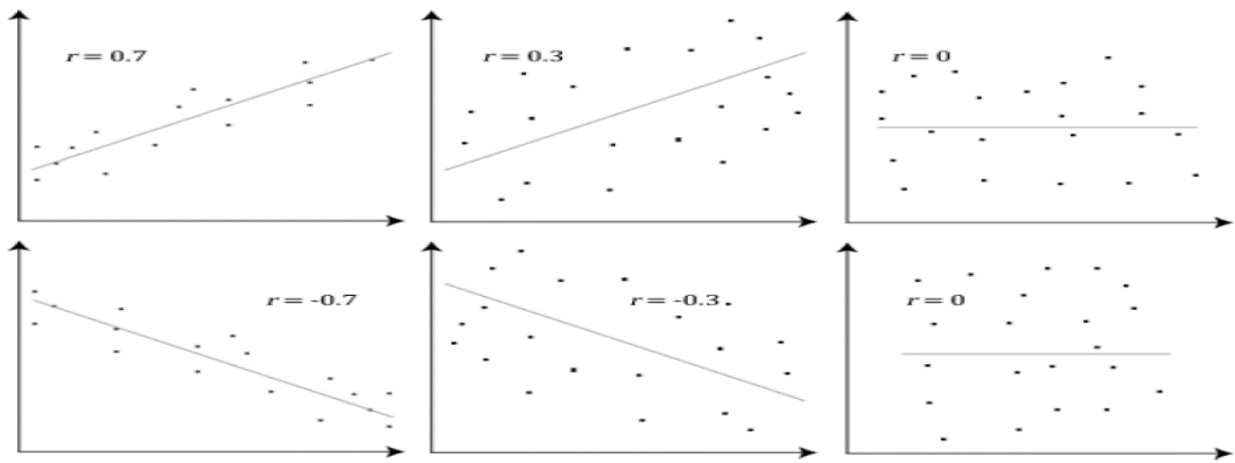
The Pearson correlation coefficient, r, can take a range of values from +1 to -1. A value of 0 indicates that there is no association between the two variables. A value greater than 0 indicates a positive association; that is, as the value of one variable increases, so does the value of the other variable.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

How can we determine the strength of association based on the Pearson correlation coefficient?

The stronger the association of the two variables, the closer the Pearson correlation coefficient, r, will be to either +1 or -1 depending on whether the relationship is positive or negative, respectively. Achieving a value of +1 or -1 means that all your data points are included on the line of best fit – there are no data points that show any variation away from this line. Values for r between +1 and -1 (for example,  $r = 0.8$  or  $-0.4$ ) indicate that there is variation around the line of best fit. The closer the value of r to 0 the greater the variation around the line of best fit. Different relationships and their correlation coefficients are shown in the diagram below:







Are there guidelines to interpreting the Pearson's correlation coefficient?

Yes, the following guidelines have been proposed: Coefficient, r

Strength of Association	Positive	Negative
Small	.1 to .3	-0.1 to -0.3
Medium	.3 to .5	-0.3 to -0.5
Large	.5 to 1.0	-0.5 to -1.0

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

**What is scaling:** 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 is scaling performed:** 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 ( $\mu$ ) zero and standard deviation one ( $\sigma$ ).

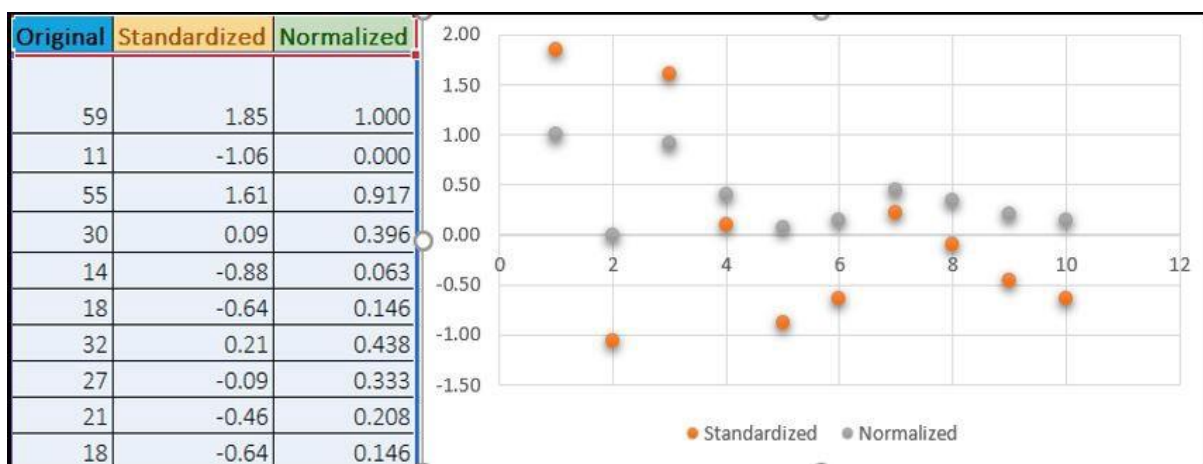
$$\text{Standardisation: } x = \frac{x - \text{mean}(x)}{\text{sd}(x)}$$

`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.

Example:

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



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

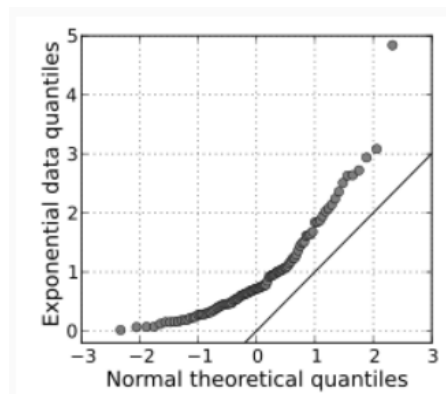
If there is perfect correlation, then  $VIF = \text{infinity}$ . This shows a perfect correlation between two independent variables. In the case of perfect correlation, we get  $R^2 = 1$ , which lead to  $1/(1-R^2)$  infinity. To solve this problem, we need to drop one of the variables from the dataset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

**6. What is a Q-Q plot? Explain the use and importance of a Q-Q plot in linear regression. (3 marks)**

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.

\*\*\*\*\*