# Bike Sharing Assignment

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

#### Answer:

- 1. The demand for bikes is very less in spring compared to other seasons.
- 2. The demand for bike is increased in 2019 compared to 2018.
- 3. Not much difference in demand between working day and non-working day.
- 4. The demand is good from June to October of the year.
- 5. Demand is little better on non-holiday than compared to holiday.
- 6. Similar for days of a week, working & non-working (not much significant differences)
- 7. Good demand in Clear or Few clouds and nearly good for Mist and Cloudy, dull on Light Snow and Rain
- 2. Why is it important to use drop\_first=True during dummy variable creation? (2 mark) Answer:

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.

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

# Answer:

Temperature(temp, also atemp) has highest correlation with 0.63 in heat map.

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

## Answer:

• Linear Relationship between the features and target:

- Little or no Multicollinearity between the features: (Bi lateral Analysis in 19th slide and heatmap in 20th slide in jupyter notebook)
- Homoscedasticity Assumption:
- Normal distribution of error terms: (Final model in 50th slide)
- Little or No autocorrelation in the residuals: You can find he correlation factors described in the final slide 52.
- 5. Based on the final model, which are the top 3 features contributing significantly towards explaining the demand of the shared bikes? (2 marks)

### Answer:

Based on final model top three features contributing significantly towards explaining the demand are:

- 1. temp (0.570)
- 2. weathersit Light Snow & Rain (-0.232)
- 3. yr (0.235)

# **General Subjective Questions**

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

## Answer:

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, the linear regression equation as

$$y = a + bx$$

Where a and b given by the formulas:

$$b(slobe) = \frac{n\sum xy - (\sum x)(\sum y)}{n\sum x^2 - (\sum x)^2}$$

$$a(intercept) = \frac{n\sum y - b(\sum x)}{n}$$

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

# 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.

Regular example used is

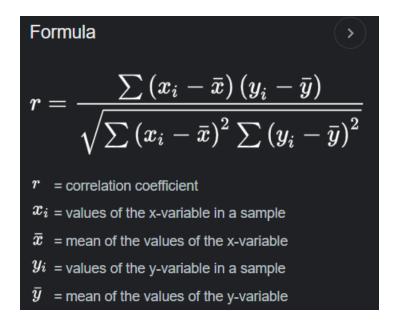
Once Francis John "Frank" Anscombe who was a statistician of great repute found 4 sets of 11 data-points in his dream and requested the council as his last wish to plot those points. Those 4 sets of 11 data-points are given below.

	I		1	II			1	III			IV		
X	1	У		X +-	7.0	У		X +-	1.00	У	1	X +-	ј у
10.0	1	8.04	1	10.0	1	9.14	1	10.0	1	7.46	1	8.0	6.5
8.0		6.95		8.0		8.14		8.0	1	6.77	1	8.0	5.7
13.0	-	7.58	1	13.0	1	8.74	1	13.0	1	12.74	1	8.0	7.7
9.0	1	8.81		9.0		8.77		9.0	-	7.11	-	8.0	8.8
11.0	1	8.33	1	11.0	1	9.26	1	11.0	1	7.81	1	8.0	8.4
14.0		9.96		14.0		8.10		14.0	1	8.84	1	8.0	1 7.0
6.0	- 1	7.24	1	6.0	1	6.13	1	6.0	1	6.08	1	8.0	5.2
4.0	1	4.26		4.0		3.10	ĺ	4.0		5.39	-	19.0	112.5
12.0	1	10.84	1	12.0	1	9.13	1	12.0	1	8.15	1	8.0	5.5
7.0	-	4.82	-	7.0		7.26		7.0	1	6.42	1	8.0	7.9
5.0		5.68	1	5.0	1	4.74	1	5.0	1	5.73	1	8.0	1 6.8

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

### Answer:

The Pearson's Correlation Coefficient is also referred to as Pearson's r, the Pearson product-moment correlation coefficient (PPMCC), or bivariate correlation. It is a statistic that measures the linear correlation between two variables. Like all correlations, it also has a numerical value that lies between -1.0 and +1.0.



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.

# 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.

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 ( $\sigma$ ).

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 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.

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

### Answer:

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.

Importance of Q-Q plot in linear regression:

- 1. It can be used with sample sizes also.
- 2. 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.
- 3. It is used to check following scenarios: (for two data sets)
  - come from populations with a common distribution
  - have common location and scale
  - have similar distributional shapes
  - have similar tail behaviour