

## 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)

From the box plot analysis

- **yr**

There are more bikes rented in 2019 as compared to 2018. This could be because of the increased popularity of the company.

- **mnth**

Clearly Jun, Jul, Aug, Sep are the months bikes are preferred more

- **season**

Fall has the highest rentals amongst all

- **weathersit**

When the weather is clear there is a trend of more rentals

- **holiday**

As bikes are more used for commute to work it seems hence the over all trend is lesser however there are chances that casuals would be using it more during holidays to move around

- **working day**

Working day has little of higher median

- **weekday**

Monday, Tuesday and Thursdays are more preferred to use bikes

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

- drop\_first=True is important to use, because it helps to reduce multicollinearity and reduces the correlations created among dummy variables. Also to avoid redundant features
- 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)

- atemp variable

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

- checked residual error which follows Normal Distribution

- also to evaluate the model checked linear relation between dependant variables with  $y_{pred}$  vs  $y_{test}$

5. Based on the final model, which are the top 3 features contributing significantly towards explaining the demand of the shared bikes? (2 marks)

- atemp with coeff of 0.4465

- yr with coeff of 0.2352

- Light rain\_Light snow\_Thunderstorm with coeff of -0.2826

## General Subjective Questions

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

- Linear Regression is a basic form of machine learning algorithm based on supervised learning. It performs a regression task. we train a model to predict the behaviour of your data based on some variables. Regression models a target prediction value based on independent variables. It is mostly used for finding out the relationship between variables and forecasting. Different regression models differ based on – the kind of relationship between dependent and independent variables they are considering, and the number of independent variables getting used.

Linear regression performs the task to predict a dependent variable value ( $y$ ) based on a given independent variable ( $x$ ). So, this regression technique finds out a linear relationship between  $x$  (input) and  $y$ (output). Hence, the name is Linear Regression.

Mathematically, we can write a linear regression equation as:

$$y = a + bx$$

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

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

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

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

## 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.

A group of four data sets which are nearly identical in simple descriptive statistics, however due to some particularities in the dataset that fools the regression model if built.

Constructed 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.

Anscombe's quartet 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.

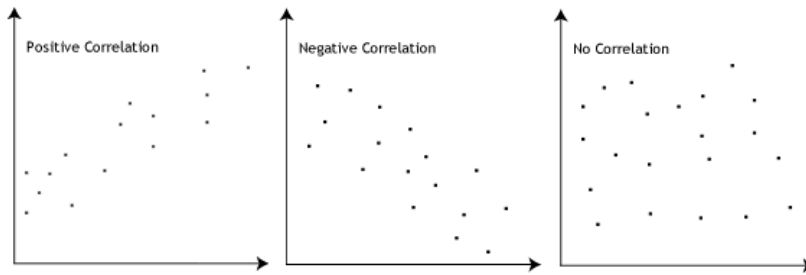
## 3. What is Pearson's R? (3 marks)

Pearson's  $r$  is a numerical summary of the strength of the linear association between the variables

In statistics, the Pearson correlation coefficient (PCC), also referred to as Pearson's  $r$ , the Pearson product-moment correlation coefficient (PPMCC), or the bivariate correlation, is a measure of linear correlation between two sets of data. It is the covariance of two variables, divided by the product of their standard deviations; thus it is essentially a normalised measurement of the covariance, such that the result always has a value between  $-1$  and  $1$ .

The Pearson's correlation coefficient varies between  $-1$  and  $+1$  where:

- $r = 1$  means the data is perfectly linear with a positive slope ( i.e., both variables tend to change in the same direction)
- $r = -1$  means the data is perfectly linear with a negative slope ( i.e., both variables tend to change in different directions)
- $r = 0$  means there is no linear association
- $r > 0 < 5$  means there is a weak association
- $r > 5 < 8$  means there is a moderate association
- $r > 8$  means there is a strong association



Pearson r Formula

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

Here,

- $r$  = correlation coefficient
- $x_i$  = values of the x-variable in a sample
- $\bar{x}$  = mean of the values of the x-variable
- $y_i$  = values of the y-variable in a sample
- $\bar{y}$  = mean of the values of the y-variable

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

- 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.

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

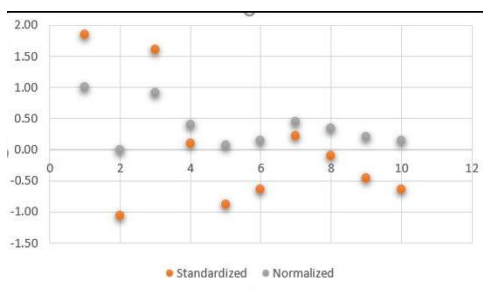
Standardization Scaling:

- 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.



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

In the case of perfect correlation, the value of VIF is infinity. As it shows a perfect correlation between two independent variables. In the case of perfect correlation, we get  $R^2 = 1$ , which leads 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)

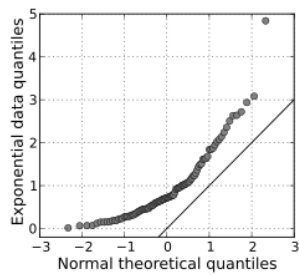
In statistics, a Q-Q plot (quantile-quantile plot) is a probability plot,

Graphical method for comparing two probability distributions by plotting their quantiles against each other.

It helps 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.

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.