

①

## Assignment-based subjective questions.

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

A) The demand of bike is less in the month of spring when compared with other seasons.

The demand bike increased in the year 2019 when compared with year 2018

Month Jun to Sep is the period when bike demand is high

The month Jan is the lowest demand month.

Bike demand is less in holidays in comparison to not being holiday.

The bike demand of bike is almost similar throughout the weekdays.

There is no significant change in the bike demand with working day and non working day.

The bike demand is high is high when weather is clear and Few clouds however demand is less in case of lightning and light rainfall.

We don't have data for Heavy Rain +~~the~~ pollution + Thunderstorm + mist + snow + fog, so we can not derive any conclusion. May be company is not operating on those days or there is no demand of bike.

- ② Why is it important to use drop-first = True during dummy variable creation?
- A) drop-first = True is important to use, as it helps in reducing the extra column created during dummy variable creation.
- ③ Looking at the pair-plot among the numerical variables which one has the highest correlation with the target variables?
- A) By looking at the pair plot temp variable has 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?
- ① Linearity
  - ② mean of Residuals
  - ③ check for Homoscedasticity
  - ④ check for Normality of error terms/residuals
  - ⑤ no autocorrelation of residuals
  - ⑥ No perfect multicollinearity
  - ⑦ other models for comparison.
- ⑤ Based on the final model, which are the top 3 features contributing significantly towards explaining the demand of the share bike?
- A) The top 3 features contributing significantly towards the demands of the share bike are:

Weather-light-Snow (negative correlation),  
yr-2019 (positive correlation).  
temp C (positive correlation).

### General subjective Questions

① Explain the linear regression algorithm in detail.

④ Linear regression is one of the very basic forms of machine learning where we train a model to predict the behavior 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

$$y = ax + b$$

where "a" and "b" given by the formula:

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

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

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 details.
- Anscombe's quartet consists of four ~~different~~ sets of data that have nearly identical simple statistical properties, yet appear very different when graphed. Each ~~dataset~~ .. Each dataset consists of eleven  $(x, y)$  points. They were constructed in 1973 by statistician Francis Anscombe to demonstrate both the importance of graphing data before analyzing it and the effect of outliers on statistical properties.

single underlying

once franc Jhen 'Frank' Anscombe who was a Statistician of great reput ~~had~~ found 4 sets of 11 ~~data~~ points in his drama and required the Council as his last wish to plot ~~those~~ those points. Those 4 sets of 11- data-points are given below.

I	II	III	IV
<u>x</u>	<u>x</u>	<u>x</u>	<u>x</u>
10.0	8.04	10.0	8.08
8.0	6.95	8.0	5.76
13.0	7.58	13.0	7.21
9.0	8.81	9.0	8.84
11.0	8.33	11.0	8.43
14.0	9.96	14.0	7.04
6.0	7.24	6.0	5.25
4.0	4.26	4.0	12.50
12.0	10.84	12.0	5.56
7.0	4.82	7.0	8.0
5.0	5.68	5.0	7.91
			8.0
			6.89

(3)

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

### ③ What is Pearson's R?

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

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

~~the~~  $r=1$  means the data is perfectly linear with a positive slope.

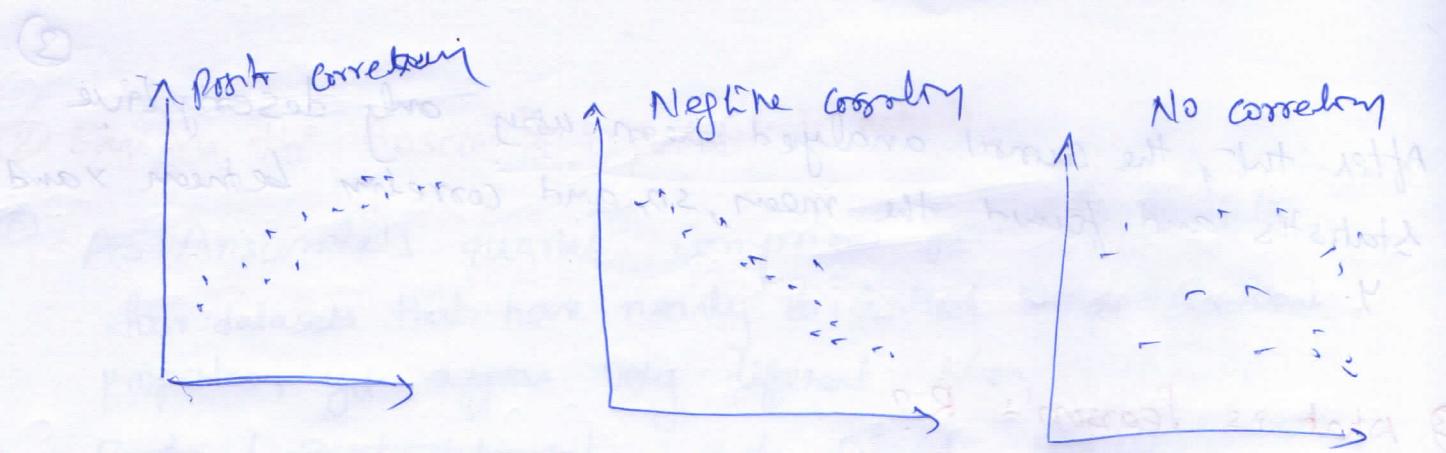
$r=-1$  means the data is perfectly linear with a negative slope.

$r=0$  means there is no linear association.

$r > 0.25$  means there is weak association

$r > 0.5 < 0.8$  means there is a moderate association

$r > 0.8$  means there is a strong association.



Pearson Formula

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

$\rho = r$  = Correlation Coefficient

$x_i$  = Value of x variable in sample

$\bar{x}$  = mean of the value of the x-variable

$y_i$  = value of the y-variable in sample

$\bar{y}$  = mean of the value of the y-variable

(4)

Q) What is the difference between normalized scaling and standardized scaling?

A) Normalized Normalization typically ~~means~~ to mean that it rescale the values into the range of  $[0, 1]$

Standardization typically means rescale the data to have a mean mean of "0" and a standard deviation of 1 (unit variance).

S. no.	Normalizing	Standardizing
①	Minimum and maximum value of feature are used for scaling.	Mean and standard deviation is used for scaling.
②	It is used when feature are in different scale.	It is used when we want to produce zeros mean and unit standard deviation.
③	Scales between $[0, 1]$ or $[-1, 1]$	It is not bounded a certain range.
④	It is really effected by outliers.	It is much less effected by outliers.
⑤	scikit-learn provides a transform called MinMaxScaler for normalizing.	scikit-learn provides a transform called StandardScaler for standardizing.
⑥	This transform squishes the <del>n-dimensional</del> data into an n-dimensional unit hypercube.	It translates the <del>data</del> to the mean vector of original <del>data</del> to the origin and squish or expand.

⑨ ~~It is useful~~

⑩ It is useful when we don't know about distribution

⑪ It is often called as Scaling Normalization

~~It is useful~~. It is useful when ~~it is~~ the feature distribution is Normal or Gaussian.

⑫ It is often called a Z-score normalization.

⑬ you ~~might~~ have observed that sometimes the value of VIF is infinity. Why does this happen?

a) If there is a perfect correlation, then  $VIF = \infty$ .

This shows a perfect correlation between two independent variables. In the case of perfect correlation, we get  $R^2 = 1$ , which lead to  $\frac{1}{(1-R^2)} = \infty$ .

In this problem we need to drop one of the variables from the data set which is causing the perfect multicollinearity.

An infinite VIF value indicates that the corresponding variable may be expressed exactly by a linear combination of other variables (which shows an infinite VIF as well).

Combination of other variables (which shows an infinite VIF as well),

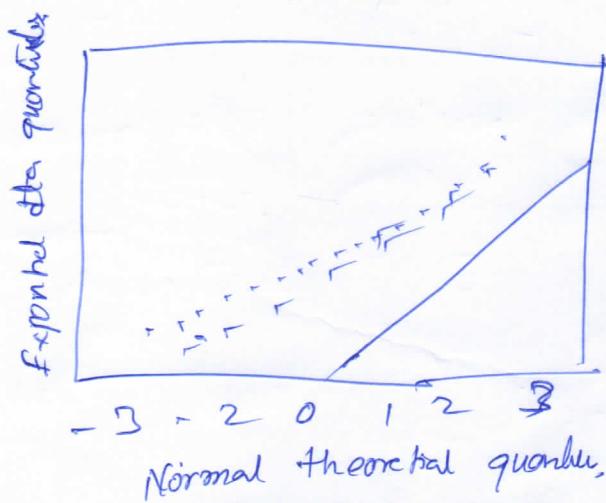
and VIF as well).

①

⑥ What is a Q-Q plot? Explain use and importance of a Q-Q plot in linear regression?

a) Q-Q plots (Quantile-Quantile plots) are plots of two quantiles against each other. A quantile is a fraction where certain values fall ~~below~~ below that quantile. For example, the median is a quantile where 50% of the ~~set~~ data fall below the point, and 50% lie above it. The purpose of Q-Q plots is to find out if two sets of data comes from the ~~same~~ data sets come from a common distribution, the points will fall on their reference line.

A Q-Q plot showing the  $45^\circ$  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 graphic means of estimating parameters in a location-scale family distributions.

A-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.

