

## 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)

**Answer:**

Below are a few points we can infer from the analysis on categorical columns:

- Most of the bookings were made during the months of May, June, July, August, September, and October. The trend exhibited an increase from the beginning of the year until mid-year, gradually decreasing as we approached the end of the year.
- The fall season appears to have attracted more bookings, while the spring season experienced fewer bookings. Additionally, the booking count has significantly increased from 2018 to 2019.
- Clear weather attracted more booking which seems obvious.
- When it's not holiday, booking seems to be less in number which seems reasonable as on holidays, people may want to spend time at home and enjoy with family.
- Booking seemed to be almost equal either on working day or non-working day.
- 2019 attracted a greater number of bookings from the previous year, which shows good progress in terms of business.

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 multicollinearity among dummy variables. This parameter is used when creating dummy variables to avoid the "dummy variable trap," which occurs when two or more dummy variables are highly correlated.

The "drop\_first" parameter, when set to True, excludes the first category, and for each categorical variable, only k-1 dummy variables are created instead of k. This helps in preventing perfect multicollinearity, where one dummy variable can be predicted perfectly from the others.

For example, if we have a categorical column with three types of values, A, B, and C, and we create dummy variables without dropping the first one, we might end up with redundant information. If a data point has values 0 for both A and B, it automatically implies the value for C without the need for a separate dummy variable. By using "`drop_first = True`," we avoid this redundancy and improve the interpretability of the model.

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

**Answer:**

'temp' variable has the highest correlation with the target variable.

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

**Answer:**

I have validated the assumptions of the Linear Regression Model based on the following five key assumptions:

- Normality of Error Terms:
  - The error terms should follow a normal distribution.
- Multicollinearity Check:
  - There should be insignificant multicollinearity among the independent variables.
- Linear Relationship Validation:
  - A linear relationship should be evident among the variables.
- Homoscedasticity:
  - There should be no discernible pattern in the residual values, ensuring constant variance.
- Independence of Residuals:
  - There should be no autocorrelation among the residuals.

By systematically validating these assumptions, I aimed to ensure that the Linear Regression model is well-suited for the dataset and provides reliable predictions on unseen data.

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:**

Below are the top 3 features contributing significantly towards explaining the demand of the shared bikes –

- temp
- year
- Light\_snowrain

## General Subjective Questions

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

**Answer:**

Linear regression may be defined as the statistical model that analyses the linear relationship between a dependent variable with given set of independent variables. The primary goal of linear regression is to find the best-fitting linear relationship that predicts the value of the dependent variable based on the values of the independent variables. Let's break down the

key components and steps involved in linear regression:

Mathematical Representation:

The linear relationship between the dependent variable (Y) and one independent variable (X) is represented by the equation:

$$Y = mX + c$$

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

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

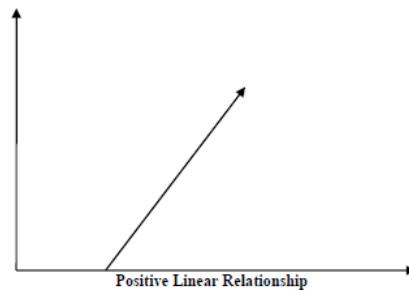
m is the slope of the regression line which represents the effect X has on Y

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

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

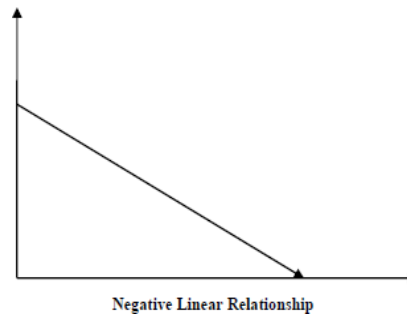
- Positive Linear Relationship:

- 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 negative if independent increases and dependent variable decreases. It can be understood with the help of following graph –



Linear regression is of the following two types –

- Simple Linear Regression
- Multiple Linear Regression

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 –

- o Linear regression model assumes that the relationship between response and feature variables must be linear.
- Normality of error terms –
  - o Error terms should be normally distributed
- Homoscedasticity –
  - o There should be no visible pattern in residual values.

## 2. Explain the Anscombe's quartet in detail.

(3 marks)

**Answer:**

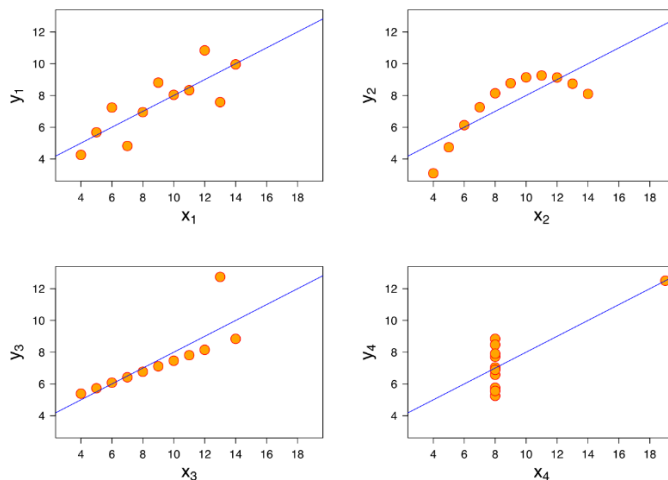
Anscombe's Quartet was developed by statistician Francis Anscombe. It comprises four datasets, each containing eleven (x, y) pairs. The essential thing to note about these datasets is that they share the same descriptive statistics. But things change completely, and I must emphasize COMPLETELY, when they are graphed. Each graph tells a different story irrespective of their similar summary statistics.

	I			II			III			IV		
	x	y		x	y		x	y		x	y	
	10	8,04		10	9,14		10	7,46		8	6,58	
	8	6,95		8	8,14		8	6,77		8	5,76	
	13	7,58		13	8,74		13	12,74		8	7,71	
	9	8,81		9	8,77		9	7,11		8	8,84	
	11	8,33		11	9,26		11	7,81		8	8,47	
	14	9,96		14	8,1		14	8,84		8	7,04	
	6	7,24		6	6,13		6	6,08		8	5,25	
	4	4,26		4	3,1		4	5,39		19	12,5	
	12	10,84		12	9,13		12	8,15		8	5,56	
	7	4,82		7	7,26		7	6,42		8	7,91	
	5	5,68		5	4,74		5	5,73		8	6,89	
SUM	99,00	82,51		99,00	82,51		99,00	82,50		99,00	82,51	
AVG	9,00	7,50		9,00	7,50		9,00	7,50		9,00	7,50	
STDEV	3,32	2,03		3,32	2,03		3,32	2,03		3,32	2,03	

The summary statistics show that the means and the variances were identical for x and y across the groups:

- Mean of x is 9 and mean of y is 7.50 for each dataset.
- Similarly, the variance of x is 11 and variance of y is 4.13 for each dataset
- The correlation coefficient (how strong a relationship is between two variables) between x and y is 0.816 for each dataset

When we plot these four datasets on an x/y coordinate plane, we can observe that they show the same regression lines as well but each dataset is telling a different story:



- Dataset I appear to have clean and well-fitting linear models.
- Dataset II is not distributed normally.
- In Dataset III the distribution is linear, but the calculated regression is thrown off by an outlier.
- Dataset IV shows that one outlier is enough to produce a high correlation coefficient.

This quartet emphasizes the importance of visualization in Data Analysis. Looking at the data reveals a lot of the structure and a clear picture of the dataset.

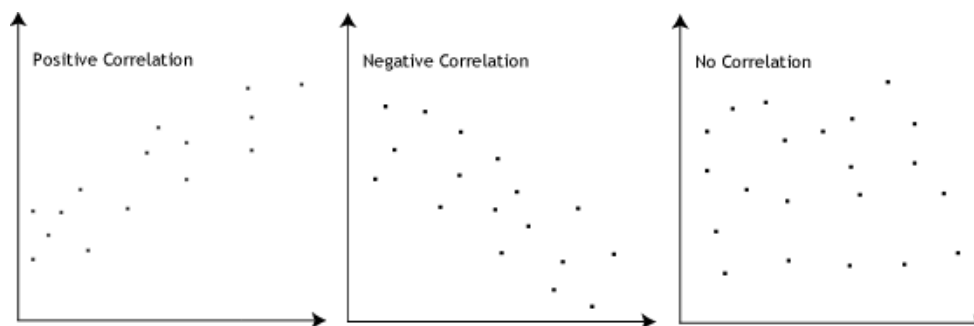
### 3. What is Pearson's R?

(3 marks)

**Answer:**

Pearson's  $r$  is a numerical summary of the strength of the linear association between the variables. If the variables tend to go up and down together, the correlation coefficient will be positive. If the variables tend to go up and down in opposition with low values of one variable associated with high values of the other, the correlation coefficient will be negative.

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:



### 4. What is scaling? Why is scaling performed? What is the difference between normalized scaling and standardized scaling?

(3 marks)

**Answer:**

Feature Scaling is a technique to standardize the independent features present in the data in a fixed range. It is performed during the data pre-processing to handle highly varying magnitudes or values or units. If feature scaling is not done, then a machine learning algorithm tends to weigh greater values, higher and consider smaller values as the lower values, regardless of the unit of the values.

Example: If an algorithm is not using feature scaling method, then it can consider the value 3000 meter to be greater than 5 km but that's actually not true and, in this case, the algorithm will give

wrong predictions. So, we use Feature Scaling to bring all values to same magnitudes and thus, tackle this issue.

S.NO.	Normalized scaling	Standardized scaling
1.	Minimum and maximum value of features are used for scaling	Mean and standard deviation is used for scaling.
2.	It is used when features are of different scales.	It is used when we want to ensure zero mean and unit standard deviation.
3.	Scales values between [0, 1] or [-1, 1].	It is not bounded to a certain range.
4.	It is really affected by outliers.	It is much less affected by outliers.
5.	Scikit-Learn provides a transformer called MinMaxScaler for Normalization.	Scikit-Learn provides a transformer called StandardScaler for standardization.

**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  $VIF = \infty$ . A large value of VIF indicates that there is a correlation between the variables. If the VIF is 4, this means that the variance of the model coefficient is inflated by a factor of 4 due to the presence of multicollinearity.

When the value of VIF is infinite 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, we need to drop one of the variables from the dataset which is causing this perfect multicollinearity.

**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 plot is used to determine whether two samples of data belong to the same population. The quantiles of the first dataset are plotted against the quantiles of the second dataset, if the two samples belong to the same population then the points will lie along the same line.

Uses of a Q-Q plot:

- Identify whether two samples belong to the same population
- Determine the distribution of the sample (normal, uniform etc...)

A q-q plot is a plot of the quantiles of the first data set against the quantiles of the second dataset. By a quantile, we mean the fraction (or percent) of points below the given value. That is, the 0.3 (or 30%) quantile is the point at which 30% percent of the data fall below and 70% fall above that value. A 45-degree reference line is also plotted. If the two sets come from a population with the same distribution, the points should fall approximately along this reference line. The greater the departure from this reference line, the greater the evidence

for the conclusion that the two data sets have come from populations with different distributions.

An example of a Q-Q plot is shown in below figure-

