# 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?

**Answer:**

Season:The average count was minimum during spring.

Year: The demand for bikes has increased in 2019 compared to 2018.

Month: The demand peaked in the month of july and then gradually decreased.

Weathersit: The demand was highest when the weather was Clear, Few clouds, Partly cloudy, Partly cloudy

**Questions 2.**

Why is it important to use drop_first=True during dummy variable creation?
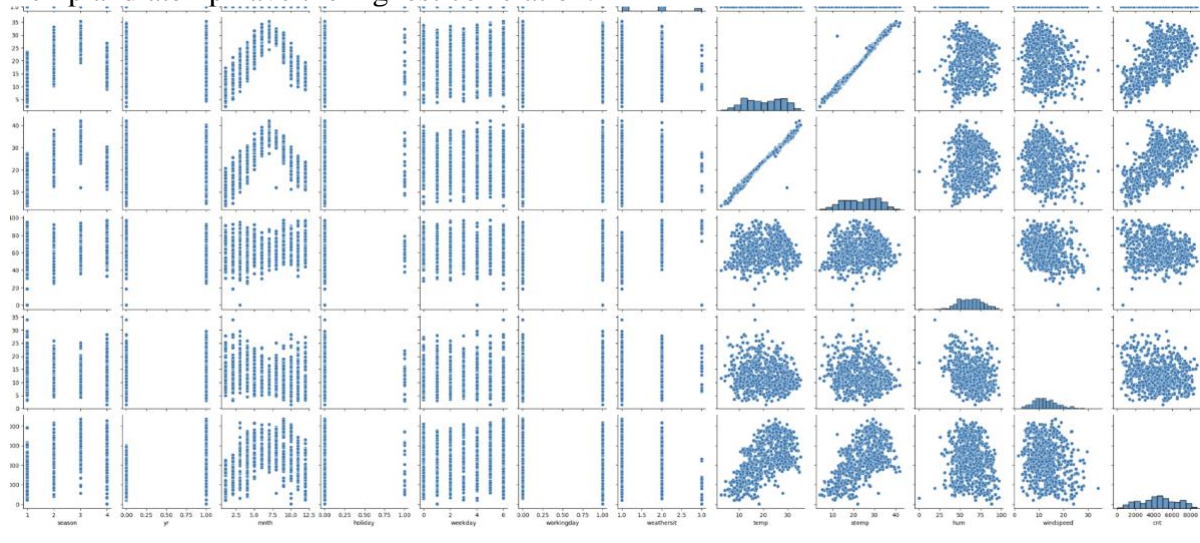
**Answer:**

For n number of variables we only need n-1 values. It helps drop extra column that gets created while creating dummy variables which might later affect the accuracy of the data.

**Questions 3.**

Looking at the pair-plot among the numerical variables, which one has the highest correlation with the target variable?

**Answer:**

Temp and atemp have the highest correlation.



**Questions 4.**

How did you validate the assumptions of Linear Regression after building the model on the training set?

**Answer:**

- Linear Relationship: The core premise of multiple linear regression is the existence of a linear relationship between the dependent (outcome) variable and the independent variables.

- No Multicollinearity: It is essential that the independent variables are not too highly correlated with each other, a condition known as multicollinearity.

- Errors are normally distributed
- Homoscedasticity of errors (or, equal variance around the line).

This can be achieved by checking vif and p values. Checking the linearity of the variables with the target variable.Getting error plot and checking if it is normally distributed.

**Questions 5.**

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

**Answer:**

1.Humidity

2. Temperature

3.Working day

# General Subjective Questions

**1. Explain the linear regression algorithm in detail**.

Linear regression is a supervised machine learning method that describes the correlation of the explanatory variables with the dependent variable. This is achieved by fitting a line to the data using least squares. The line tries to minimize the sum of the squares of the residuals. The residual is the distance between the line and the actual value of the explanatory variable. Finding the line of best fit is an iterative process.

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

Anscombe's quartet is a group of datasets (x, y) that have the same mean, standard deviation, and regression line, but which are qualitatively different.

The datasets were created by the statistician Francis Anscombe in 1973 to demonstrate the importance of visualizing data and to show that summary statistics alone can be misleading.

The four datasets that make up Anscombe's quartet each include 11 x-y pairs of data. When plotted, each dataset seems to have a unique connection between x and y, with unique variability patterns and distinctive correlation strengths. Despite these variations, each dataset has the same summary statistics, such as the same x and y mean and variance, x and y correlation coefficient, and linear regression line.

Purpose of Anscombe's Quartet:

Anscombe's quartet is used to illustrate the importance of exploratory data analysis and the drawbacks of depending only on summary statistics. It also emphasizes the importance of using data visualization to spot trends, outliers, and other crucial details that might not be obvious from summary statistics alone.

Anscombe's Quartet Dataset:

The four datasets of Anscombe's quartet.:

```
+-------+--------+-------+-------+-------+--------+-------+------+
|     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 |
+-------+--------+-------+-------+-------+--------+-------+------+
```

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

The Pearson correlation coefficient (r) is the most common way of measuring a linear correlation. It is a number between –1 and 1 that measures the strength and direction of the relationship between two variables. When one variable changes, the other variable changes in the same direction.

The Pearson correlation coefficient ($r$) is the most widely used correlation coefficient and is known by many names:

- Pearson's $r$
- Bivariate correlation
- Pearson product-moment correlation coefficient (PPMCC)
- The correlation coefficient

Although interpretations of the relationship strength (also known as efffect size) vary between disciplines, the table below gives general rules of thumb:

| Pearson correlation coefficient ($r$) value | Strength | Direction |
|---|---|---|
| Greater than .5 | Strong | Positive |
| Between .3 and .5 | Moderate | Positive |
| Between 0 and .3 | Weak | Positive |

| Pearson correlation coefficient (*r*) value | Strength | Direction |
| --- | --- | --- |
| 0 | None | None |
| Between 0 and –.3 | Weak | Negative |
| Between –.3 and –.5 | Moderate | Negative |
| Less than –.5 | Strong | Negative |

The Pearson correlation coefficient is also an inferential statistic meaning that it can be used to test statistical hypothesis. Specifically, we can test whether there is a significant relationship between two variables.

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

Scaling is the process of expanding or compressing the dimensions of an object (changing the size of an object).Scaling is performed to negate different scales used in the dataset.

Normalization and standardization are two techniques used to transform data into a common scale. Normalization is a technique used to scale numerical data in the range of 0 to 1. This technique is useful when the distribution of the data is not known or when the data is not normally distributed. On the other hand, standardization is a technique used to transform data into a standard normal distribution. This technique is useful when the distribution of the data is known and when the data is normally distributed. Both techniques have different applications, and choosing the right technique based on the data and the problem you're trying to solve is important.

It is a scaling method that reduces duplication in which the numbers are scaled and moved between 0 and 1. When there are no outliers since it can't handle them, normalization is employed to remove the undesirable characteristics from the dataset.

Normalization, occasionally is a method for rescaling the values that meet the characteristics of the standard normal distribution while being similar to normalizing.

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

If there is perfect correlation, then VIF = infinity. A large value of VIF indicates that there is a correlation between the 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.**

Q–Q plots are commonly used to compare a data set to a theoretical model. This can provide an assessment of goodness of fit that is graphical, rather than reducing to a numerical summary statistic. Q–Q plots are also used to compare two theoretical distributions to each other.

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.
This helps in a scenario of linear regression when we have training and test data set received separately and then we can confirm using Q-Q plot that both the data sets are from populations with same distributions.

A q-q plot is a plot of the quantiles of the first data set against the quantiles of the second data set.