# Assignment-based Subjective Questions

1. <mark>From your analysis of the categorical variables from the dataset, what could you infer about their effect on the dependent variable? (3 marks)</mark>

Analysing categorical variables in a dataset can provide insights into how these variables might influence the dependent variable. Typically, this is done by examining the relationship between the categories and the dependent variable. Here's how you can approach the analysis and what you might infer:

**Steps to Analyse Categorical Variables:**

1. **Encoding Categorical Variables**:
   o Convert categorical variables into numerical formats (e.g., one-hot encoding, label encoding).
   o Use drop first=True in one-hot encoding to avoid the dummy variable trap (perfect multicollinearity).
2. **Summary Statistics**:
   o Calculate min & max values of the dependent variable across different categories.
   o Analyse how the average values differ between categories.
3. **Visualization**:
   o Create boxplots to visualize the distribution of the dependent variable across different categories.
   o Look for patterns, such as certain categories consistently having higher or lower values of the dependent variable.
4. **Feature Importance**:
   o After fitting the model, analyse the coefficients of the categorical variables (in linear regression) or feature importance (in tree-based models) to understand their impact.

**Categorical Values Box Plot Analysis**

**1**. Season:

   • Fall having higher rental range where spring having less rental range.

2. Year:

   • 2019 Rental of bike considerably increased compare to 2018

3. Weather:

   • Rental having higher value during Clear weather than cloudy mist and light rain thunder.

4. WeekDay:

   • Saturday and Thursday more rental than other working days.

5. HoliDay:

- - Holiday rented less than working days.

6. Month:

- September leads in monthly rental counts, with nearby months also showing strong demand. This trend follows seasonal patterns, suggesting a link between rentals and seasons.

2. <mark>Why is it important to use drop_first=True during dummy variable creation? (2 mark)</mark>

Using `drop_first=True` during dummy variable creation is important to avoid the **dummy variable trap**, which can lead to **multicollinearity** in the model.

- **Avoiding Multicollinearity:**

  When creating dummy variables, if all categories are included, the sum of the dummy variables can perfectly predict one another, leading to multicollinearity. This can cause instability in the model coefficients and affect the interpretability of the results.

- **Ensuring a Reference Category:**

  By dropping the first dummy variable, one category is used as a reference, and the remaining categories are compared against this reference. This ensures that the model remains well-posed, with one less predictor, improving model interpretability and reducing redundancy.

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

Highest correlation with cnt is seen in temp.

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

❑ Linearity:

- **Assumption**:
  - The relationship between the independent variables and the dependent variable is linear.
- **Validation**:
  - **Plot Residuals vs. Fitted Values**: This plot should show no clear pattern. If the residuals are randomly scattered around zero, it suggests that the linearity assumption holds.
  - **Scatter Plots**: Plot each independent variable against the dependent variable to visually check for linear relationships.

❑ Normality :

- **Assumption:**

- o  The residuals are normally distributed.
- **Validation:**

Histogram of Residuals: The histogram of the residuals should look like a bell curve (normal distribution).

❑ Multicollinearity :

- **Assumption:**
    - o  The independent variables are not too highly correlated with each other. Variable must be independent of each other
- **Validation:**
    - o  Variance Inflation Factor (VIF): Calculate the VIF for each predictor. A VIF value greater than 10 indicates high multicollinearity and suggests that the predictor may need to be removed or combined with others

❑ Homoscedasticity :

- **Assumption:**
    - o  **The variance of the residuals is constant across all levels of the independent** variables.
- **Validation:**
    - o  Plot Residuals vs. Fitted Values: In addition to checking linearity, this plot is also used to check homoscedasticity. The spread of residuals should be constant across all levels of fitted values.

❑ Autocorrelation :

- o  Autocorrelation is the correlation between two of the same series. Residual Analysis of the train 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)

Our final multiple linear regression model seeks to predict bike bookings using a collection of predictor factors.

Following a thorough investigation, we determined the influence of each variable on bike bookings.

The top three variables(with co-efficient) that significantly influence bike rent are as follows:

- Yr 2036.64
- Temp 106.00
- Windspeed -35.00

# General Subjective Questions

Linear regression is a supervised learning algorithm used for predicting a continuous outcome based on one or more input features (independent variables). The goal is to establish a linear relationship between the dependent variable (target) and the independent variables by fitting a straight line (in simple linear regression) or a hyperplane (in multiple linear regression) that best represents the data.

**The equation of a linear regression** model can be represented as:

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \cdots + \beta_n x_n + \epsilon$$

Where:

- $y$ is the predicted output (dependent variable).
- $x_1, x_2, \ldots, x_n$ are the input features (independent variables).
- $\beta_0$ is the intercept (the value of $y$ when all $x_i = 0$).
- $\beta_1, \beta_2, \ldots, \beta_n$ are the coefficients (slopes) that represent the change in $y$ for a unit change in $x_i$.
- $\epsilon$ is the error term, representing the difference between the observed and predicted values.

**Training the Model**

The linear regression model is trained by finding the values of $\beta_1, \beta_2, \ldots, \beta_n$ that minimize the cost function, typically the Mean Squared Error (MSE):

$$MSE = \frac{1}{m} \sum_{i=1}^{m} (y_i - \hat{y}_i)^2$$

Where:

- $m$ is the number of data points.
- $y_i$ is the actual value.
- $\hat{y}_i$ is the predicted value.

The minimization is usually done using optimization techniques like Gradient Descent, which iteratively updates the coefficients to reduce the error.

**Prediction and Interpretation**

Once trained, the model can predict the outcome for new data by substituting the values of the independent variables into the regression equation. The coefficients $\beta_i$ indicate the strength and direction of the relationship between each feature and the target. For example, a positive $\beta_i$ suggests that as $x_i$ increases, $y$ also increases, while a negative $\beta_i$ suggests the opposite.

In summary, linear regression is a foundational algorithm in machine learning, providing a simple yet powerful method for predicting outcomes and understanding relationships between variables.

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. It is often used to illustrate the importance of looking at a set of data graphically and not only relying on basic statistic properties.

It highlights the importance of data visualization to gain a deeper understanding of the underlying patterns and relationships within the data.

3. <mark>What is Pearson's R? (3 marks)</mark>

Pearson's R, also known as the Pearson correlation coefficient, is a statistical measure that quantifies the strength and direction of the linear relationship between two continuous variables. It ranges from -1 to +1.

- **Value Interpretation:**
  - +1: Indicates a perfect positive linear relationship, where as one variable increases, the other also increases in a perfectly linear manner.
  - -1: Indicates a perfect negative linear relationship, where as one variable increases, the other decreases in a perfectly linear manner.
  - 0: Indicates no linear relationship between the variables.
- **Importance**:
  - In regression analysis, it helps in assessing the strength of the linear relationship between the independent and dependent variables.

4. <mark>What is scaling? Why is scaling performed? What is the difference between normalized scaling and standardized scaling? (3 marks)</mark>

<mark>What is scaling</mark>

Scaling is a step of data Pre-Processing. This applied over independent variables to normalize the data within a particular range. It also helps in speeding up the calculations in an algorithm.

<mark>Why is scaling performed</mark>

In most cases, the collected data contains highly diverse characteristics in terms of scale, units and scope. If scaling is not done, the algorithm only considers the magnitude and not the units, resulting in inaccurate model. To solve this problem, we need to perform some scaling to bring all variables to the same level of magnitude.

<mark>What is the difference between normalized scaling and standardized scaling</mark>

Standardization centres data around a mean of zero and a standard deviation of one, while normalization scales data to a set range, often [0, 1], by using the minimum and maximum values. One disadvantage of normalization over standardization is that it loses some information in the data, especially about outliers.

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

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. To solve this problem we need to drop one of the variables from the dataset which is causing this perfect multicollinearity.

6. <mark>What is a Q-Q plot? Explain the use and importance of a Q-Q plot in linear regression.</mark>

<mark>(3 marks)</mark>

A Q-Q (Quantile-Quantile) plot is a graphical tool used to compare the distribution of a dataset with a theoretical distribution, typically the normal distribution. It plots the quantiles of the dataset against the quantiles of the theoretical distribution. If the data follows the theoretical distribution, the points on the Q-Q plot will approximately lie on a straight line.

**Importance in Linear Regression:** In linear regression, one of the key assumptions is that the residuals (the differences between observed and predicted values) are normally distributed. A Q-Q plot helps assess this assumption by comparing the quantiles of the residuals to the quantiles of a normal distribution. If the residuals are normally distributed, the points in the Q-Q plot will follow a straight line.