# 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?
Ans :

The categorical variables such as **season**, **yr**, **mnth**, **holiday**, **weekday**, **workingday**, and **weathersit** have varying levels of influence on bike demand. For instance:

- **Bike rentals** tend to increase during the **autumn** months.
- The months of **June, August, and October** experience a surge in **bike rentals**.
- Favorable **weather conditions** often lead to a higher preference for **renting bikes**.
- **Weekends** see a lower tendency for people to **rent bikes**.
- **Bike rentals** are less popular on **holidays**.
- The year **2019** saw **more bike rentals** compared to the year **2018**

2. Why is it important to use drop_first=True during dummy variable creation?

Ans:
a) **Avoiding Multicollinearity**:
   - When we create dummy variables for categorical features without dropping one category, it introduces **multicollinearity**.
   - By dropping one category (usually the reference category), we can prevent multicollinearity.
   - The information from the dropped category is inherently included in the remaining dummy variables.
b) **Enhancing Interpretability**:
   - Setting drop_first=True makes the interpretation of model coefficients more intuitive.
   - It ensures that the reference category is excluded, allowing us to focus on the effects of other categories relative to the reference.

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

Ans:
When we look at the pair-plot among the numerical variables, we notice that the **"registered"** variable has a **strong correlation** with the target variable **"cnt"**. It's like they're close friends who always hang out together in the scatter plot!

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

Ans:

1. **Avoiding Multicollinearity**:
   - In statistical modeling, multicollinearity occurs when two or more independent variables are highly correlated with each other. It's like having friends who always agree and say the same things—it doesn't provide new insights.
   - When we create dummy variables for categorical features (like "Neighborhood" with categories such as "Downtown," "Suburb," and "Rural"), we need to drop one category.

Why? Because if we keep all categories, they become linearly dependent, causing multicollinearity.

o By dropping one category (usually the reference category), we ensure that the information from the dropped category is inherently included in the remaining dummy variables. It's like saying, "Okay, let's focus on the unique aspects of each friend without redundancy."

2. **Enhancing Interpretability**:
o Imagine explaining your model to someone over coffee. When you set drop_first=True, you're making the conversation smoother.
o Let's say "Downtown" is our reference category. Now, when you see a positive coefficient for "Suburb," you can confidently say, "Ah, houses in the suburbs tend to have higher prices compared to downtown." It's like decoding a secret message!
o So, dropping that first category makes the model coefficients more intuitive and easier to explain. It's like translating complex math into everyday language.

Top features that significantly contribute to explaining the demand for shared bikes:

• **Temperature (temp):**

The variable 'temp' is indicative of the ambient temperature conditions. It has a substantial impact on bike rental demand, as favorable temperatures are conducive to outdoor activities such as cycling.

• **Weather Situation (weathersit_bad):**

The 'weathersit_bad' variable captures adverse weather conditions. Its significant negative influence on bike rental demand underscores the sensitivity of this activity to weather disruptions.

• **Year (yr):**

The 'yr' variable represents the year of observation. Its inclusion as a top feature suggests a temporal trend in bike rental demand, possibly reflecting an increase in the popularity of bike-sharing services over time.

These features underscore the importance of environmental and temporal factors in the utilization of bike-sharing services.

**General Subjective Questions**

1. Explain the linear regression algorithm in detail.

Ans :
Linear Regression is a statistical technique in machine learning that models the relationship between a dependent variable and one or more independent variables. It's particularly useful in projects like bike-sharing demand prediction, where you want to understand how factors such as temperature, humidity, and time of day affect the number of bikes rented.

**Simple vs. Multiple Linear Regression**

• **Simple Linear Regression** uses one predictor to forecast outcomes. For example, predicting bike rentals based on temperature alone with the equation

Bike Rentals=b0+b1·Temperature
.

• **Multiple Linear Regression** takes into account several factors, leading to a more complex equation like

Bike Rentals=b0+b1·Temperature+b2·Humidity+…+bn·Other Factors
.

**Assumptions and Model Fitting**
The model assumes a linear relationship, independence of errors, homoscedasticity (constant variance of errors), and normally distributed errors. It aims to minimize the sum of squared differences between actual and predicted values, finding the best-fitting line or hyperplane.

**Insights for Bike-Sharing Demand**
By analyzing the coefficients, we can determine which factors have the most significant impact on bike rentals. This helps in making informed decisions about inventory management and promotional strategies, ensuring that the service is optimized for user demand.
In essence, Linear Regression provides a clear, quantifiable way to predict and understand the patterns in bike-sharing demand, enabling efficient service planning and resource allocation.

2. Explain the Anscombe's quartet in detail.

Anscombe's quartet is a set of four datasets that are designed to demonstrate the importance of visualizing data before analyzing it. Despite having nearly identical statistical properties—such as mean, variance, and correlation—each dataset looks very different when graphed.

In the context of a bike-sharing demand prediction project, Anscombe's quartet serves as a cautionary tale. It reminds us that relying solely on summary statistics could lead to incorrect conclusions. For instance, if we were to analyze factors affecting bike rentals using only mean and variance, we might miss out on underlying patterns or anomalies that could be crucial for accurate predictions.

Therefore, before building a predictive model, it's essential to plot the data. This way, we can identify if the relationship between variables like temperature, humidity, or time of day and bike rentals is linear (suitable for linear regression) or if there are outliers or non-linear patterns that require a different analytical approach. Visual analysis ensures that the chosen model is appropriate for the data's actual distribution and relationships.

3. What is Pearson's R?

Ans :

Pearson's Correlation Coefficient ($r$) is a statistical measure that assesses the strength and direction of the linear relationship between two continuous variables. It ranges from -1 to 1, where:
($r = 1$) indicates a perfect positive linear correlation.
($r = -1$) indicates a perfect negative linear correlation.
($r = 0$) suggests no linear correlation.
In the context of bike-sharing demand prediction, ($r$) could be used to evaluate how strongly factors like temperature or humidity are linearly related to the number of bikes rented. A high positive ($r$) value with temperature would suggest that higher temperatures lead to more rentals, guiding decisions on inventory and marketing strategies.

- What is scaling? Why is scaling performed? What is the difference between normalized scaling and standardized scaling?

Ans :

**Scaling in Machine Learning**
- **Definition**: Scaling is a data preprocessing technique used to transform feature values to a similar scale.
- **Purpose**:
  - Ensure that all features contribute equally to the model.
  - Prevent features with larger values from dominating the learning process.
  - Facilitate meaningful comparisons between features.
  - Improve model convergence.
  - Prevent numerical instability.
- **Methods**: Two common scaling techniques are **Normalization** and **Standardization**.

**Normalization (Min-Max Scaling)**
- **Formula**: $X_{\text{new}} = \frac{X - X_{\text{min}}}{X_{\text{max}} - X_{\text{min}}}$
- **Range**: Scales values between [0, 1] or sometimes [-1, 1].
- **Use Case**:
  - When features have different scales.
  - No assumption about the underlying data distribution.
  
  **Example**:
  - Scaling bike rental counts to [0, 1].

**Standardization (Z-Score Normalization)**
- **Formula**: $X_{\text{new}} = \frac{X - \text{mean}}{\text{Std}}$
- **Range**: Not bounded to a certain range.
- **Use Case**:
  - When we want zero mean and unit standard deviation.
  - Feature distribution is assumed to be normal or Gaussian.
  
  **Example**:
  - Standardizing temperature and humidity.

In summary, scaling ensures fair feature contributions, faster model convergence, and robustness against outliers. Choose normalization for diverse features and standardization for Gaussian-distributed features

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

Ans :

The Variance Inflation Factor (VIF) assesses the level of multicollinearity among predictor variables in a regression model. A VIF value becomes infinite when there is perfect multicollinearity, meaning one predictor variable can be exactly predicted by a combination of other variables.

In the context of a bike-sharing demand prediction project, if we have a VIF that is infinite, it could be due to a situation where one of the features, such as the number of bikes available, is an exact linear combination of other features, like the number of docks and the number of bikes already rented out. This perfect predictability indicates redundant data, which can distort the regression results and make the coefficients unreliable.

To address this, we would need to identify and remove the redundant variables, ensuring that each feature provides unique information about the bike-sharing demand. This step is crucial for developing a robust predictive model.

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

Ans:

A Q-Q plot, short for "quantile-quantile" plot, is a graphical tool used to assess whether a set of data potentially came from some theoretical distribution, typically a normal distribution1. In linear regression, it's used to check the normality of residuals2.

In the bike-sharing demand prediction project, a Q-Q plot can be crucial for validating the assumption that the residuals (differences between observed and predicted rental counts) are normally distributed. If the points in the Q-Q plot lie on or near a straight diagonal line, it suggests that the residuals are normally distributed1. This normality is important because many inferential statistics and hypothesis tests in linear regression rely on this assumption2.

If the residuals are not normally distributed, it could indicate that the linear model may not be the best fit for the data, and the reliability of the regression results could be compromised. Therefore, a Q-Q plot is an essential diagnostic tool in ensuring the robustness of the bike-sharing demand prediction model12.