

## Assignment-based Subjective Questions

**Question 1.** From your analysis of the categorical variables from the dataset, what could you infer about their effect on the dependent variable? (Do not edit)

**Total Marks:** 3 marks (Do not edit)

**Answer:** <Your answer for Question 1 goes below this line> (Do not edit)

Based on the categorical variables in the dataset ('season', 'yr', 'mnth', 'holiday', 'weekday', 'workingday', and 'weathersit'), below are the effects on the dependent variable, 'cnt' (bike demand):

1. Season:

- Inference: Demand is typically lower in spring and winter and higher in summer and fall.

2. Year:

- Inference: The increase in demand from 2018 to 2019 likely indicates growing acceptance of bike-sharing, though an outlier in 2019 could indicate a demand dip due to the pandemic's onset.

3. Month:

- Inference: Warmer months (e.g., May through September) likely show increased demand, while colder months show reduced usage.

4. Holiday:

- Inference: Demand is generally lower on holidays, as fewer people commute to work, suggesting bikes are more commonly used for commuting than leisure.

5. Weekday:

- Inference: Demand might be slightly higher on weekdays, reflecting consistent commuting patterns, but the distribution across weekdays is typically similar.

6. Working Day:

Inference: Demand is usually higher on working days than non-working days, again highlighting the bike service's utility for commuting.

7. Weather Situation:

- Inference: Demand is highest in clear or mild weather (weathersit = 1) and decreases as conditions worsen, with the lowest demand for the harshest weather (weathersit = 3).

---

**Question 2.** Why is it important to use **drop\_first=True** during dummy variable creation? (Do not edit)

**Total Marks:** 2 marks (Do not edit)

**Answer:** <Your answer for Question 2 goes below this line> (Do not edit)

Using 'drop\_first=True' when creating dummy variables is important because it helps avoid the dummy variable trap, which occurs due to multicollinearity. Here's why it matters:

1. Multicollinearity Prevention:

- Creating a dummy variable for each level means that they would be highly correlated. Knowing the values of all but one level lets you perfectly determine the value of the last level.
- Including all dummy variables would lead to perfect multicollinearity in linear models, making it difficult for the model to distinguish between the effects of each dummy variable.

2. Interpretation of the Baseline:

- Dropping the first dummy variable makes one category the reference or baseline category, against which other categories are compared.
- For instance, if we drop the dummy for "spring" in 'season', the model interprets "spring" as the baseline, and coefficients for "summer", "fall", and "winter" represent their relative effect compared to spring.

In short, 'drop\_first=True' ensures our model avoids redundancy and multicollinearity, giving it a more stable foundation for interpreting categorical effects.

---

**Question 3.** Looking at the pair-plot among the numerical variables, which one has the highest correlation with the target variable? (Do not edit)

**Total Marks:** 1 mark (Do not edit)

**Answer:** <Your answer for Question 3 goes below this line> (Do not edit)

In a pair-plot among numerical variables, 'registered' usually has the highest correlation with the target variable 'cnt' (total bike rentals), as 'cnt' is the sum of 'casual' and 'registered' rentals. Here's a breakdown:

'registered' and 'cnt': 'registered' users (those with memberships or regular usage) often show a strong positive correlation with 'cnt' because they constitute a large, consistent portion of total rentals.

---

**Question 4.** How did you validate the assumptions of Linear Regression after building the model on the training set? (Do not edit)

**Total Marks:** 3 marks (Do not edit)

**Answer:** <Your answer for Question 4 goes below this line> (Do not edit)

After building the model on the training set, here's how I validated each assumption:

1. Linearity:

Method: Plot the predicted values against the residuals.

2. Homoscedasticity (Constant Variance of Errors):

Method: Examine the residual plot (residuals vs. predicted values).

3. Normality of Residuals:

Method: Plot a histogram of residuals.

4. No Multicollinearity:

Method: Calculated the Variance Inflation Factor (VIF) for each predictor.

---

**Question 5.** Based on the final model, which are the top 3 features contributing significantly towards explaining the demand of the shared bikes? (Do not edit)

**Total Marks:** 2 marks (Do not edit)

**Answer:** <Your answer for Question 5 goes below this line> (Do not edit)

Here are the top three features that contributed significantly to explaining bike demand ('cnt')

1. registered:

Explanation: This variable represents the count of registered users and generally has the strongest correlation with 'cnt', as it directly contributes to the total rentals. Registered users often provide consistent demand, which explains much of the variability in bike usage.

2. casual: casual rentals also correlate with cnt, but typically to a lesser extent than registered, since casual rentals tend to fluctuate more based on external factors (like weather and holidays)

3. 'temp':

Explanation: Temperature is a significant predictor because people are more likely to use bikes in moderate weather. Demand typically increases with favorable temperatures, making it a key factor in the final model.

---

## General Subjective Questions

**Question 6.** Explain the linear regression algorithm in detail. (Do not edit)

**Total Marks:** 4 marks (Do not edit)

**Answer:** Please write your answer below this line. (Do not edit)

<Your answer for Question 6 goes here>

Linear regression is a supervised learning algorithm used primarily for predictive analysis and modeling the relationship between a dependent variable and one or more independent variables. It assumes a linear relationship between the input variables (predictors) and the single output variable (response). Let's go through it in detail:

1. Objective of Linear Regression

The main goal of linear regression is to find the line that best fits the data points. This line (the regression line) can then be used to predict the output variable for new data points. The equation of the line in simple linear regression (one predictor) is:

$$y = b_0 + b_1x_1 + b_2x_2 + \dots$$

---

- $y$ : Predicted output (dependent variable)
- $x$ : Independent variable (input)
- $b_0, b_1, b_2$ : Coefficients of  $x$

## 2. Working of Linear Regression

Linear regression tries to estimate the best values for  $b_0, b_1, b_2, \dots$  by minimizing the error in predictions. Specifically, it minimizes the sum of squared differences (residuals) between the predicted and actual values. This method is known as Ordinary Least Squares (OLS).

**Residuals (Errors):** The difference between the actual values and the predicted values.

**Cost Function (Mean Squared Error):** The cost function represents the average squared difference between actual and predicted values and helps assess model accuracy. The objective is to minimize this MSE to find the best-fitting line.

## 3. Gradient Descent Optimization

Gradient Descent is an optimization algorithm often used to minimize the cost function. It iteratively adjusts the coefficients by calculating the gradient (slope) of the cost function and taking steps toward a minimum cost. The update rule for the parameters is:

This process continues until the cost function converges to a minimum value.

## 4. Assumptions of Linear Regression

For linear regression to be reliable, certain assumptions need to hold:

- **Linearity:** The relationship between independent and dependent variables should be linear.
- **Independence:** Observations should be independent of each other.
- **Homoscedasticity:** The variance of error terms should be constant across all levels of the independent variables.
- **Normality of Errors:** Errors should follow a normal distribution.
- **No Multicollinearity (in Multiple Linear Regression):** Independent variables should not be highly correlated with each other.

## 5. Evaluation Metrics for Linear Regression

- **R-squared:** Measures the proportion of the variance in the dependent variable that is predictable from the independent variables. It ranges from 0 to 1, with 1 indicating a perfect fit.
- **Adjusted R-squared:** Adjusts for the number of predictors in the model, penalizing models that add predictors without improving the model significantly.
- **Mean Absolute Error (MAE):** The average of absolute differences between actual and predicted values.
- **Root Mean Squared Error (RMSE):** The square root of the average squared differences between actual and predicted values. RMSE is more sensitive to outliers than MAE.

## 6. Applications of Linear Regression

- Predicting sales or stock prices based on historical data
  - Forecasting trends in time series data
-

- Evaluating the impact of variables on economic indicators
- Estimating risk in insurance and finance
- Modeling relationships in scientific data, such as dosage-response analysis in pharmacology

Linear regression is a foundational model in machine learning due to its simplicity and interpretability. It fits a line through data by minimizing prediction errors and is widely used in predictive analytics and inferential statistics. By meeting its assumptions, linear regression can yield insights and serve as a baseline for more complex models in regression analysis.

---

**Question 7.** Explain the Anscombe's quartet in detail. (Do not edit)

**Total Marks:** 3 marks (Do not edit)

**Answer:** Please write your answer below this line. (Do not edit)

<Your answer for Question 7 goes here>

Anscombe's quartet is a set of four datasets that have nearly identical summary statistics (mean, variance, correlation, and linear regression line) but display very different patterns when graphed. The quartet was created by the British statistician Francis Anscombe in 1973 to demonstrate the importance of visualizing data before analyzing it solely through statistical summaries. This concept underscores that statistical properties alone can be misleading if data isn't properly visualized and inspected.

#### Datasets in Anscombe's Quartet

Each dataset in Anscombe's quartet has:

- The same mean of the  $x$ -values.
- The same mean of the  $y$ -values.
- The same variance of both  $x$ - and  $y$ -values.
- The same correlation between  $x$  and  $y$ .
- The same linear regression line (when  $y$  is regressed on  $x$ ).

Despite these identical statistics, each dataset has a distinct visual pattern, which tells a different story about the data. Let's take a closer look at each dataset in the quartet:

#### 1. Dataset 1

- Description: This dataset resembles what would be expected from typical data suitable for linear regression, with points lying roughly along a straight line.
- Characteristics: It has a clear linear relationship with points scattered around a line.
- Key Insight: This dataset justifies the use of linear regression as a modeling approach, showing that a simple linear relationship can be effectively summarized by regression analysis.

#### 2. Dataset 2

- Description: This dataset has a nearly perfect linear relationship but with one clear outlier.
- Characteristics: The points are mostly in a vertical line (same  $x$ -values) except for one point.
- Key Insight: The presence of an outlier greatly affects the regression line and correlation, emphasizing the importance of checking for outliers before using statistical measures like

correlation or regression.

### 3. Dataset 3

- Description: This dataset has a clear curve rather than a linear relationship.
- Characteristics: The data follows a strong nonlinear, quadratic pattern.
- Key Insight: Although linear regression gives a line, it is not a good model for this dataset due to the underlying nonlinear relationship, highlighting the limitations of linear regression on nonlinear data.

### 4. Dataset 4

- Description: This dataset has most points with the same  $y$ -value except for one outlier.
- Characteristics: Nearly all the points lie on a single vertical line, with one point far away.
- Key Insight: The outlier greatly influences the regression line and correlation, masking the fact that most of the data does not have a meaningful relationship. This dataset shows the sensitivity of regression and correlation to single extreme values.

### Key Lessons from Anscombe's Quartet

1. Importance of Data Visualization: Relying on statistical summaries alone can lead to misleading interpretations. Visualizations such as scatter plots reveal patterns, outliers, and relationships that statistics alone may not.

2. Impact of Outliers: Outliers can heavily influence statistical measures, especially in small datasets, potentially skewing insights.

3. Limitations of Linear Regression: Linear regression might not be appropriate for all datasets, particularly those with nonlinear relationships. Visual inspection can help confirm if linear regression is suitable.

4. Role of Context in Data Analysis: Different data structures (linear, nonlinear, presence of outliers) can have different interpretations, so understanding the data context is crucial.

### Summary

Anscombe's quartet is a powerful reminder that numbers alone don't tell the full story, and visualizing data can uncover nuances that would otherwise be overlooked. It reinforces the importance of complementing statistical analysis with data visualization for more accurate and reliable interpretations.

---

**Question 8.** What is Pearson's  $R$ ? (Do not edit)

**Total Marks:** 3 marks (Do not edit)

**Answer:** Please write your answer below this line. (Do not edit)

<Your answer for Question 8 goes here>

Pearson's  $r$ , also known as the Pearson correlation coefficient, is a measure of the linear relationship between two variables. It ranges from -1 to 1 and quantifies both the strength and direction of this relationship. This measure is commonly used in statistics to assess how closely two variables are related to each other.

### ### Interpreting Pearson's r

The value of  $r$  indicates the type and strength of the relationship between the variables:

- $r = 1$  : Perfect positive linear relationship. As  $X$  increases,  $Y$  also increases proportionally.
- $r = -1$  : Perfect negative linear relationship. As  $X$  increases,  $Y$  decreases proportionally.
- $r = 0$  : No linear relationship between  $X$  and  $Y$ .

### Properties of Pearson's r

1. Symmetry: The correlation between  $X$  and  $Y$  is the same as between  $Y$  and  $X$ .
2. Scale Independence: The value of  $r$  remains the same if  $X$  and/or  $Y$  are scaled linearly.
3. Sensitive to Outliers: Outliers can strongly influence  $r$ , making it important to assess outliers before interpreting  $r$  values.

### ### Assumptions for Pearson's r

For Pearson's  $r$  to accurately describe the relationship, certain conditions should hold:

- Linearity: The relationship between  $X$  and  $Y$  should be linear.
- Independence: Observations should be independent of each other.
- Normality:  $X$  and  $Y$  should be approximately normally distributed (especially in smaller samples).
- Homoscedasticity: The variance of  $Y$  should be roughly the same for all values of  $X$ .

### ### Example Use Cases of Pearson's r

- Assessing the correlation between height and weight.
- Understanding the relationship between study hours and test scores.
- Evaluating the connection between temperature and electricity consumption.

### ### Summary

Pearson's  $r$  is a simple, effective measure to gauge the linear relationship between two variables. However, visualizing data alongside calculating  $r$  is essential for a complete understanding, as  $r$  does not capture non-linear relationships or account for outliers.

The formula for Pearson's correlation coefficient is:

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

where:

- $x_i$  and  $y_i$  are individual data points for variables  $X$  and  $Y$ .
- $\bar{x}$  and  $\bar{y}$  are the means of  $X$  and  $Y$ .
- The summations and square roots capture the relationship between the deviations of  $X$  and  $Y$  from their respective means.

---

**Question 9.** What is scaling? Why is scaling performed? What is the difference between normalized scaling and standardized scaling? (Do not edit)

**Total Marks:** 3 marks (Do not edit)

**Answer:** Please write your answer below this line. (Do not edit)

<Your answer for Question 9 goes here>

Scaling is a preprocessing technique used in machine learning and data science to adjust the range of data features to ensure consistency across them. It's especially important when working with algorithms that rely on distance measures (e.g., k-NN, SVM) or those sensitive to feature variance (e.g., gradient descent-based models).

Scaling helps in:

1. Improving Model Performance: Many machine learning models perform better when data is scaled because they rely on distances or feature variance. Scaling ensures no feature dominates due to its scale.
2. Accelerating Convergence: In optimization algorithms, such as gradient descent, scaling helps achieve faster convergence because features of similar scales allow the algorithm to move more smoothly through the parameter space.
3. Preventing Bias: Models are less likely to be biased toward features with larger scales, which can otherwise skew predictions.

Types of Scaling: Normalization vs. Standardization

There are two common methods of scaling data: Normalization and Standardization. Here's a breakdown of each and their key differences.

#### 1. Normalization (Min-Max Scaling)

Normalization, often called min-max scaling, transforms data to fit within a specific range, typically [0, 1] or [-1, 1].

$$X' = \frac{X - \mu}{\sigma}$$

where:

- $X'$  is the standardized value.
- $X$  is the original data point.
- $\mu$  is the mean of the feature.
- $\sigma$  is the standard deviation of the feature.

Characteristics:

- Range-Bounded: Normalization transforms values to a specified range, making it particularly useful for algorithms sensitive to the absolute range, like neural networks or k-NN.
- Sensitive to Outliers: Outliers can distort the min-max range, causing most data to be scaled into a narrow range.

Use Cases:

- When feature ranges need to be preserved.
- When data doesn't contain extreme outliers.



## 2. Standardization (Z-Score Scaling)

Standardization, also called Z-score scaling, transforms data so that it has a mean of 0 and a standard deviation of 1.

$$X' = \frac{X - X_{\min}}{X_{\max} - X_{\min}}$$

where:

- $X'$  is the normalized value.
- $X$  is the original data point.
- $X_{\min}$  and  $X_{\max}$  are the minimum and maximum values of the feature.

Characteristics:

- Zero-Centered: Standardized data has a mean of 0 and unit variance, which helps center data, making it compatible with many machine learning algorithms, especially those assuming normally distributed data.
- Less Affected by Outliers: While not completely immune, standardization is generally less sensitive to outliers than normalization.

Use Cases:

- Suitable when features have very different ranges or when data has outliers.
- Often preferred in algorithms like linear regression, logistic regression, and SVMs, where normally distributed data yields better results.

### Summary of Differences

Feature	Normalization	Standardization
Range	Specific range, typically [0, 1] or [-1, 1]	Mean 0, Standard Deviation 1
Use Cases	Distance-based models, neural networks	Normal distribution-based models
Effect of Outliers	Sensitive to outliers	Less sensitive to outliers
Focus	Rescales values relative to range	Centers around mean and scales by variance

When to Use Which?

- Normalization is generally used when the features are bounded and distance-based methods are used.
- Standardization is more common when data distribution is approximately normal or if there are outliers present.

**Question 10.** You might have observed that sometimes the value of VIF is infinite. Why does this happen? (Do not edit)

**Total Marks:** 3 marks (Do not edit)

**Answer:** Please write your answer below this line. (Do not edit)

<Your answer for Question 10 goes here>

An infinite Variance Inflation Factor (VIF) value typically indicates perfect multicollinearity among the predictors in a regression model. VIF is a measure of how much the variance of a regression coefficient is inflated due to multicollinearity in the model.

Understanding VIF and Multicollinearity

The VIF for a predictor is calculated as:

$$VIF = 1 / (1 - R^2)$$

where  $R^2$  is the coefficient of determination (R-squared)

- When  $R^2_i = 1$ , meaning the predictor  $X_i$  is perfectly linearly predictable from the other predictors, the denominator becomes zero, making the VIF value for  $X_i$  infinite.
- This situation arises when there is perfect multicollinearity—in other words, one predictor is an exact linear combination of the other predictors.

Why Infinite VIF Occurs

Infinite VIF values occur due to:

1. Redundant Predictors: When two or more predictors are linearly dependent or when one variable is an exact multiple of another (e.g., height in inches and height in centimeters).
2. Over-parameterization: Too many predictors relative to the number of observations or when some predictors do not add new information to the model, creating perfect linear dependence.

Consequences and Solution

Infinite VIF values signal that the regression model cannot reliably estimate the effect of the collinear predictors. This issue affects model stability, interpretability, and predictive accuracy.

Solutions:

- Remove redundant predictors: Identify and remove perfectly collinear predictors to reduce multicollinearity.
- Dimensionality reduction techniques: Methods like Principal Component Analysis (PCA) can help in reducing multicollinearity by transforming correlated predictors into uncorrelated components.

Detecting and handling multicollinearity ensures a more reliable and interpretable model by mitigating inflated VIF values.

---

**Question 11.** What is a Q-Q plot? Explain the use and importance of a Q-Q plot in linear regression. (Do not edit)

**Total Marks:** 3 marks (Do not edit)

**Answer:** Please write your answer below this line. (Do not edit)

<Your answer for Question 11 goes here>

A Q-Q (Quantile-Quantile) plot is a graphical tool used to assess if a dataset follows a particular theoretical distribution, most commonly the normal distribution. In a Q-Q plot, the quantiles of the dataset are plotted against the quantiles of the theoretical distribution. If the points in the Q-Q plot lie approximately on a straight line, it indicates that the data follows the specified distribution.

#### How a Q-Q Plot Works

1. **Quantiles of Data:** The data is ordered, and quantiles (specific data points corresponding to specific cumulative probabilities) are calculated.
2. **Theoretical Quantiles:** Quantiles from the theoretical distribution (e.g., normal distribution) are computed based on the same probabilities as the data quantiles.
3. **Plotting:** The sample quantiles (from the data) are plotted on the y-axis against the theoretical quantiles on the x-axis.

If the data aligns well with the theoretical distribution, the points should follow a roughly straight line.

#### Importance of Q-Q Plots in Linear Regression

In linear regression, the assumption that residuals (errors) are normally distributed is crucial for accurate interpretation and prediction. A Q-Q plot can help verify if this assumption holds:

1. **Checking Normality of Residuals:** The normality of residuals is assumed in linear regression to ensure reliable hypothesis testing, confidence intervals, and prediction intervals. If the Q-Q plot shows significant deviations from a straight line, it suggests that residuals are not normally distributed, which can impact the validity of p-values and confidence intervals.
2. **Identifying Skewness or Heavy Tails:** Q-Q plots reveal if the distribution is skewed (curved line), or if it has heavy tails or light tails (points deviating from the line at the extremes). These characteristics indicate that the model might require transformations or a different approach, such as robust regression.
3. **Detecting Outliers:** Deviations from the straight line, especially at the ends, can signal outliers in the data. Outliers can have a significant impact on linear regression results, affecting coefficient estimates and predictions.

#### How to Interpret a Q-Q Plot in Linear Regression

- **Straight Line:** Residuals are approximately normally distributed, satisfying the normality assumption.
- **S-Shaped Curve:** Indicates skewness. Right skewness if the plot curves downward to the left and upward to the right, and left skewness if the curve is in the opposite direction.
- **Heavy Tails:** If the plot flares out at the ends, it indicates heavy tails (more extreme values than a normal distribution).
- **Light Tails:** If the plot compresses at the ends, it suggests light tails (fewer extreme values than a normal distribution).

#### Summary

A Q-Q plot is a valuable diagnostic tool in linear regression, allowing analysts to check for

normality and potential outliers in residuals. Ensuring normally distributed residuals through Q-Q plot analysis leads to more robust and interpretable regression models.

---