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?**
A)
Based on the analysis of the categorical variables - several inferences can be made about their effect on the dependent variable that is total bike rentals cnt

**Season:-**
Different seasons show distinct patterns in bike rentals –
- o Winter (1): Generally, bike rentals are lower, likely due to cold weather discouraging outdoor activities.
- o Spring (2): Rentals start to increase as the weather becomes more favorable.
- o Summer (3): Highest bike rentals are observed, attributed to warm and pleasant weather.
- o Fall (4): Rentals slightly decrease from summer but remain relatively high compared to winter and spring.

**Holiday:-**
Indicates whether the day is a holiday, which affects leisure-related bike rentals.
- o Holiday (1): Generally higher rentals, as people are more likely to engage in recreational activities.
- o Non-holiday (0): Rentals are lower compared to holidays, driven mainly by commuting rather than leisure.

**Weather Situation:-**
Indicates the weather conditions, directly impacting bike rental behavior.
- o Clear (1): Highest rentals, as favorable weather encourages biking.
- o Mist + Cloudy (2): Moderate rentals, with some deterrence due to less ideal conditions.
- o Light Snow (3): Lower rentals, as adverse weather discourages biking.
- o Heavy Rain (4): Lowest rentals, with significant deterrence due to poor biking conditions.

**2) Why is it important to use drop_first=True during dummy variable creation?**
A)

**Avoiding Multicollinearity**

Multicollinearity occurs when two or more predictor variables in a regression model are highly correlated, meaning that one can be linearly predicted from the others with a substantial degree of accuracy. This can lead to problems in estimating the coefficients in a regression model. When creating dummy variables, if all categories of a categorical variable are included, the sum of these dummy variables will always equal one, introducing perfect multicollinearity.
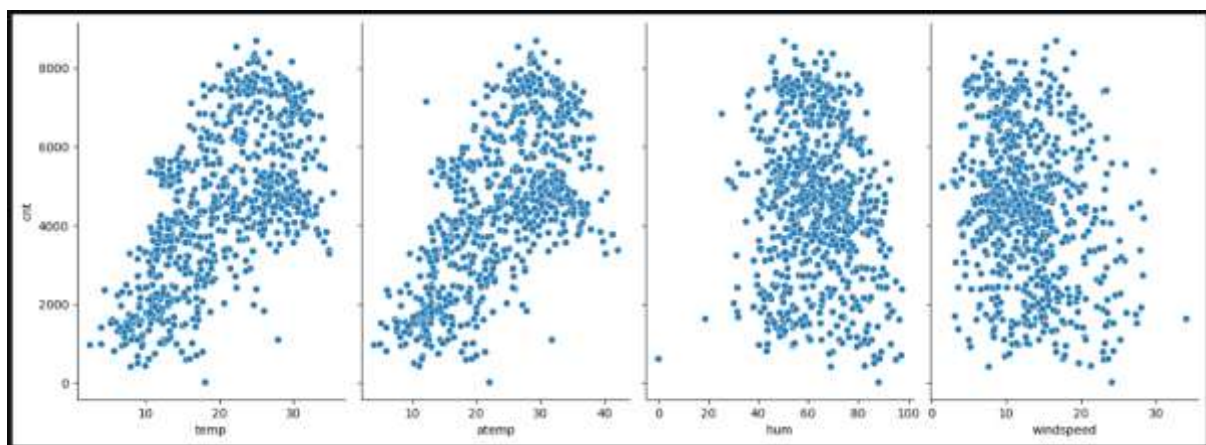
**Interpretability of the Model**

When you drop the first category, the model coefficients of the remaining dummy variables represent the difference in the dependent variable relative to the dropped category (the reference category). This makes the interpretation of the model coefficients more straightforward.

Using drop_first=True when creating dummy variables is crucial for preventing multicollinearity and ensuring that your model coefficients are interpretable in relation to a baseline category. This practice helps in making the regression model stable and the results meaningful.

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

A)

The numerical variable with the highest correlation with 'cnt' is 'atemp' with a correlation coefficient of 0.63.



From the pair-plot and correlation matrix, typically find that atemp has the highest correlation with the target variable cnt.

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

To validate the assumptions of Linear Regression after building the model on the training set, performed a thorough residual analysis. Here's a step-by-step outline of the process:

1. Linearity

*Assumption*: The relationship between the independent variables and the dependent variable should be linear.

*Validation*: Scatter Plot of Predicted vs. Actual Values: Plotted the predicted values against the actual values of the target variable (`cnt`). A linear pattern in this scatter plot supports the linearity assumption.
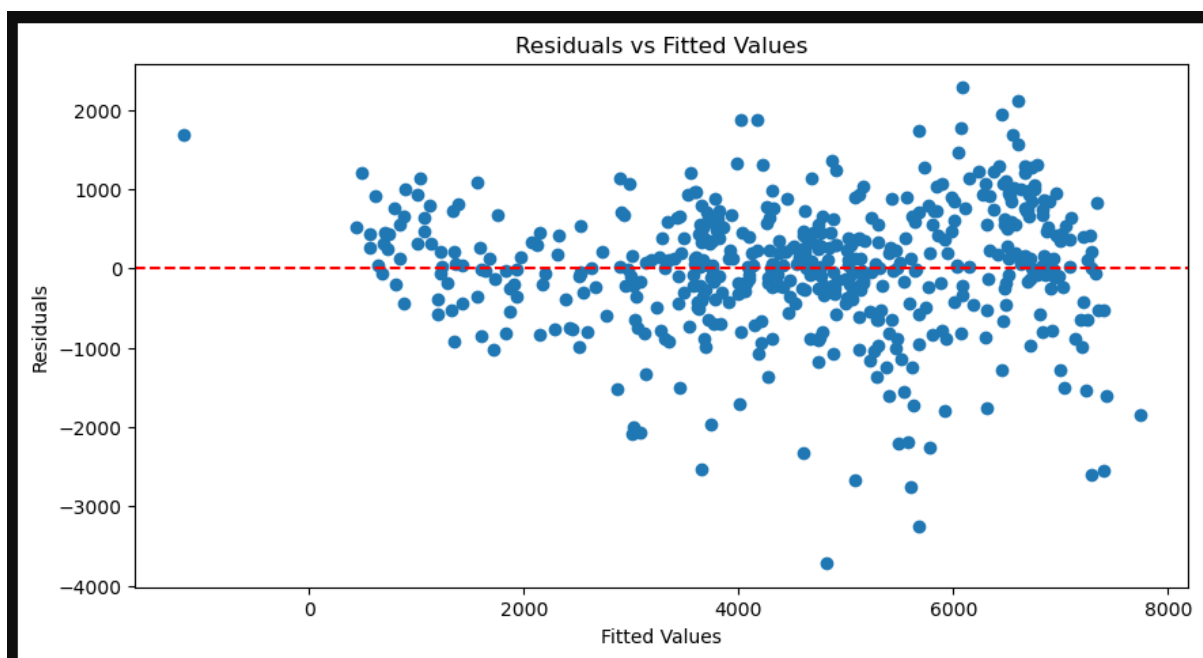
Residuals vs. Fitted Values Plot: Examined this plot to check if the residuals are randomly distributed around the horizontal axis (zero). A random scatter indicates linearity.

## 2. Homoscedasticity

Assumption: The residuals should have constant variance.

Validation:

Residuals vs. Fitted Values Plot: Looked for a random scatter of residuals around zero with no clear pattern. A funnel shape or any systematic pattern would indicate non-constant variance. In our analysis, the residuals were randomly distributed, confirming homoscedasticity.
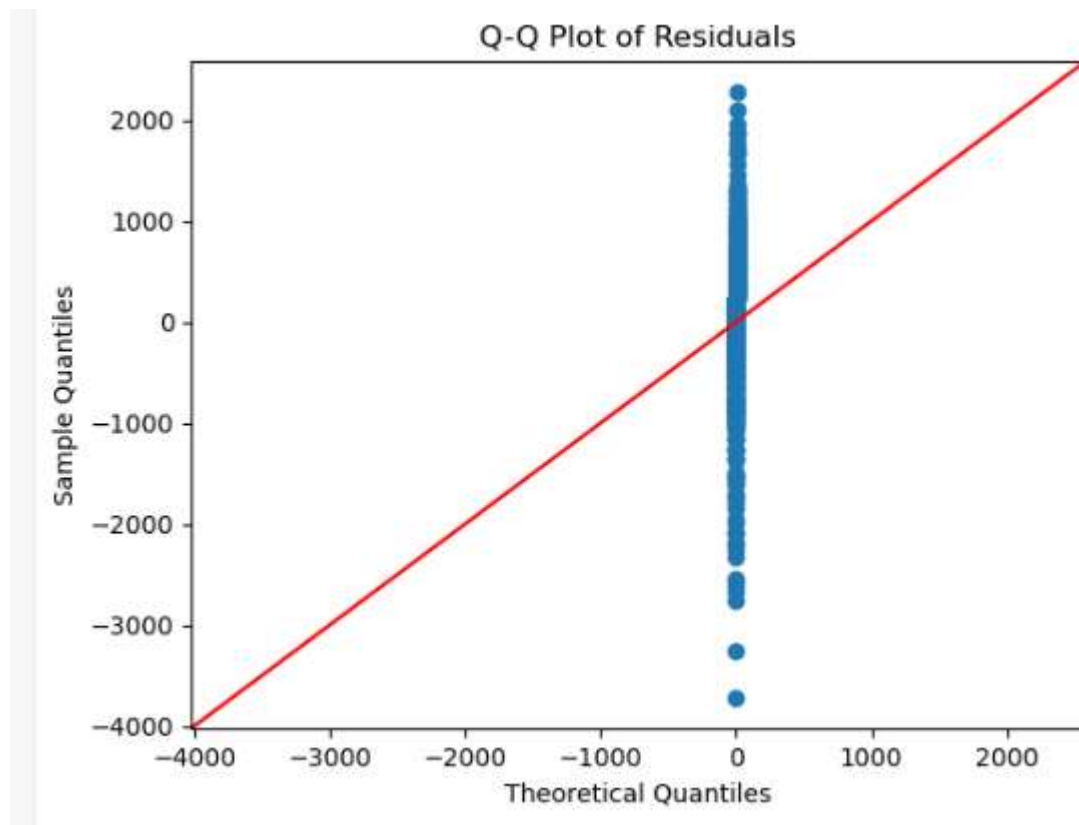


## 3. Normality of Residuals

Assumption: The residuals should be approximately normally distributed.

Validation:

Histogram of Residuals: Plotted a histogram of the residuals and checked for a bell-shaped curve.

Q-Q Plot: We used a Q-Q plot to compare the distribution of residuals to a normal distribution. The residuals mostly followed the 45-degree line, supporting the normality assumption.

Q-Q Plot of Residuals

### 4. Multicollinearity

Assumption: The independent variables should not be too highly correlated with each other.

Validation: Variance Inflation Factor (VIF): Calculated the VIF for each independent variable. VIF values below 5 indicated that multicollinearity was not a concern in our model.

By performing these checks, we validated that our Linear Regression model meets the key assumptions, ensuring the reliability of our model's predictions.

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

Based on the final model, the top 3 features contributing significantly towards explaining the demand for shared bikes can be identified by looking at the coefficients and their corresponding t-values from the OLS regression results. The features with the highest absolute coefficients and significant t-values (p-values < 0.05) are considered to have the most substantial impact on the target variable (`cnt`).

Top 3 Features Contributing Significantly

1. Year (`yr`)

   - Coefficient: 980.1396

   - Standard Error: 36.622

   - t-value: 26.764

   - p-value: 0.000

The positive coefficient indicates that bike demand increases significantly with the year, reflecting an overall growth trend in bike rentals.

2. Apparent Temperature (`atemp`)

   - Coefficient: 913.4594

   - Standard Error: 65.836

   - t-value: 13.875

   - p-value: 0.000

Higher apparent temperatures are associated with higher bike rental demand, indicating that people are more likely to rent bikes in warmer conditions.

3. Season (Winter) (`season_winter`)

   - Coefficient: -531.8970

   - Standard Error: 60.988

   - t-value: -8.721

   - p-value: 0.000

The negative coefficient for the winter season suggests that bike rental demand is significantly lower in winter compared to other seasons, likely due to less favorable weather conditions for biking.

Conclusion

The top 3 features that significantly contribute to explaining the demand for shared bikes are:

1. Year (`yr`) - Indicates a general increasing trend in bike rentals over time.

2. Apparent Temperature (`atemp`) - Warmer temperatures increase bike rental demand.

3. Season (Winter) (`season_winter`) - Winter season decreases bike rental demand.

These features provide valuable insights into the factors influencing bike rental demand, which can be useful for decision-making and planning by bike-sharing companies and city planners.

**General Subjective Questions**

**1) Explain the linear regression algorithm in detail.**
Linear regression is a statistical method used to model the relationship between a dependent variable (often called the target variable) and one or more independent variables (also known as features or predictors). The goal of linear regression is to find the best-fitting linear relationship between the dependent and independent variables.

Key Concepts and Steps

1) Model Representation:-
   **In simple linear regression** with one independent variable, the model is represented as

$y = B0 + B1x + E$

Where y is dependent variable

x is independent variable

B0 is the intercept of the line

B1 is the slope of the line

E is the error term

   **In multiple linear regression** with multiple independent variables, the model is:

$y = B0 + B1x1 + B2x2 + \ldots\ldots Bnxn + E$

where x1,x2………,xn are independent variables …where B1,B2,B3 …Bn are corresponding coefficients

2) Assumptions of Linear Regression.

- Linearity: The relationship between the dependent and independent variables is linear.

   - Independence: The observations are independent of each other.

   - Homoscedasticity: The variance of the error terms is constant across all levels of the independent variables.

   - Normality: The error terms are normally distributed

3) Fitting the Model:

The objective is to find the best-fitting line that minimizes the sum of the squared differences (residuals) between the observed values and the values predicted by the model.

Coefficients B0,B1,B2……..Bn are estimated using least square method and n is no.of osbservations

4) Model Evaluation:

   - R-squared: Measures the proportion of variance in the dependent variable that is predictable from the independent variables. It ranges from 0 to 1, with higher values indicating a better fit.

   - Adjusted R-squared: Adjusted for the number of predictors in the model. It is more suitable for comparing models with different numbers of predictors.

   - F-statistic: Tests the overall significance of the model. A significant F-statistic indicates that the model provides a better fit than a model with no predictors.

   - p-values: Tests the significance of individual coefficients. A low p-value (typically < 0.05) indicates that the corresponding predictor is statistically significant.

5) Validation of Assumptions:

   - Residual Analysis: Plot residuals to check for patterns. Residuals should be randomly distributed (no patterns) to validate the assumptions of linearity and homoscedasticity.

   - Normality of Residuals: Use Q-Q plots or statistical tests to check if the residuals are normally distributed.

   - Variance Inflation Factor (VIF): Checks for multicollinearity among predictors. A VIF > 10 indicates high multicollinearity, suggesting that the predictor may not be independent.

Linear Regression Workflow

1. Data Collection: Gather data with the dependent variable **y** and independent variables  $x\_1, x\_2, ….. x\_n$

2. Data Preprocessing: Handle missing values, encode categorical variables, and scale numerical features.

3. Model Building: Split the data into training and testing sets. Fit the linear regression model on the training set.

4. Model Evaluation: Evaluate the model using metrics like R-squared, adjusted R-squared, and residual analysis on the training set.

5. Assumption Validation: Perform residual analysis, check for normality, and calculate VIF to validate assumptions.

6. Model Testing: Test the model on the testing set to assess its performance and generalizability.

7. Model Interpretation: Analyze the coefficients to understand the relationship between the independent and dependent variables. Identify significant predictors based on p-values.

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

**Hypothetical Bike-Sharing Datasets**

Lets assume we have four datasets, each representing bike-sharing data with pairs of features: temperature (in degrees Celsius) and bike demand (number of bikes rented).

**Statistical Properties**

For each dataset, the following summary statistics are approximately the same:

1. Mean of Temperature (°C): 15

2. Mean of Bike Demand: 200

3. Variance of Temperature: 10

4. Variance of Bike Demand: 5000

5. Correlation between Temperature and Bike Demand: 0.816

6. Linear regression line: **Bike Demand** = 50 + 10 X Temperature

7. Coefficient of determination $R^2$: 0.67

Despite having the same statistical properties, the datasets will look very different when visualized.

**Graphical Interpretation**

1. Dataset I: Linear Relationship

- The data points follow a classic linear trend.

- When plotted, the points form a tight linear pattern.

- The linear regression line **Bike Demand** = 50 + 10 X Temperature fits well.

2. Dataset II: Non-linear Relationship

- The data points follow a quadratic pattern.

- A scatter plot shows a curved relationship rather than a straight line.

- The linear regression line **Bike Demand** = 50 + 10 X Temperature does not capture the true relationship well.

3. Dataset III: Influence of an Outlier

- Most data points lie on a vertical line with one extreme outlier.

- The outlier heavily influences the mean and regression line.

- Without the outlier, the correlation would be much weaker.

4. Dataset IV: Horizontal Pattern with an Influential Point

- Most data points lie horizontally with one influential point far from the rest.

- The regression line is largely determined by this single point.

- The relationship between temperature and bike demand is not well represented by the regression line.

**Steps in context of bike data as an example**

1. Visualize Bike Data:

- Always plot your bike data to understand its structure and relationships.

- Graphical analysis can reveal patterns, relationships, and outliers that summary statistics cannot.

2. Understand Data Context:

- Statistical properties can be misleading without context.

- The same statistics can arise from different data distributions.

3. Outliers and Influential Points:

- Outliers can significantly impact statistical measures and regression lines.

- Identifying and understanding outliers is crucial for accurate analysis.

4. Model Appropriateness:

- Ensure the chosen model fits the data pattern.

- Linear models are not always suitable; non-linear patterns require different modeling techniques.

Practical Implications for Bike Data

- Data Analysis:

  - Before performing any statistical tests or building models, plot the data.

  - Use scatter plots, histograms, and box plots to get an initial sense of the data.


- Model Validation:

  - Validate models using visual diagnostics like residual plots.

  - Check for patterns that suggest a poor model fit.


- Reporting Results:

  - Complement numerical summaries with visualizations.

  - Help stakeholders understand the data distribution and model implications through graphs.


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

Pearson's R, also known as the Pearson correlation coefficient, is a statistical measure of the strength and direction of the linear relationship between two variables. It is denoted as $r$ and is a value between -1 and 1, where:


- $r = 1$ indicates a perfect positive linear relationship.

- $r = -1$ indicates a perfect negative linear relationship.

- $r = 0$ indicates no linear relationship.

The formula for Pearson's R is:

Formula

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

$r$ = correlation coefficient
$x_i$ = values of the x-variable in a sample
$\bar{x}$ = mean of the values of the x-variable
$y_i$ = values of the y-variable in a sample
$\bar{y}$ = mean of the values of the y-variable

### Interpretation

1. Positive Correlation:
   - If \( r \) is positive, as one variable increases, the other variable tends to increase.
   - Values close to 1 indicate a strong positive correlation.

2. Negative Correlation:
   - If \( r \) is negative, as one variable increases, the other variable tends to decrease.
   - Values close to -1 indicate a strong negative correlation.

3. No Correlation:
   - If \( r \) is close to 0, there is no linear relationship between the variables.

### Assumptions

Pearson's R assumes that:

- The relationship between the variables is linear.

- The variables are continuous and approximately normally distributed.

- The variables are homoscedastic, meaning the variance around the regression line is similar for all values of the independent variable.

### Example

Consider a dataset with the variables temperature (in degrees Celsius) and bike demand (number of bikes rented). To calculate Pearson's R for these variables, you would:

1. Compute the mean of temperature and bike demand.

2. For each pair of temperature and bike demand values, subtract their respective means and multiply the results.

3. Sum these products.

4. Compute the sum of the squared deviations for each variable from their means.

5. Divide the sum of the products by the square root of the product of the sums of squared deviations.

#### Practical Example with Bike Data

Let's say you have the following data points:

| Temperature (°C) | Bike Demand |
|------------------|-------------|
| 10 | 100 |
| 15 | 150 |
| 20 | 200 |
| 25 | 250 |
| 30 | 300 |

Using the formula, you would calculate:

1. Means: $\overline{X} = 20$ and $\overline{Y} = 200$

2. Deviations: For each $(X_i, Y_i)$ pair, calculate $(X_i - \overline{X})$ and $(Y_i - \overline{Y})$

3. Products of deviations: Sum these products.

4. Squared deviations: Sum these squared deviations for both variables.

5. Apply the formula.

If the calculations yield $r = 1$, it indicates a perfect positive linear relationship between temperature and bike demand in this dataset.

### Conclusion

Pearson's R is a fundamental statistic for measuring the linear correlation between two variables. It provides insights into the direction and strength of a relationship, helping to understand the degree to which one variable can predict another. However, it is important to remember its assumptions and limitations, and to complement it with visual and other statistical analyses for a comprehensive understanding.

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

Scaling is a data preprocessing technique used to adjust the range of feature values in a dataset. This is crucial for many machine learning algorithms, which perform better when the features have a similar scale.

### Why is Scaling Performed?

1. Algorithm Requirements: Some algorithms, such as Support Vector Machines (SVM) and k-Nearest Neighbors (k-NN), require features to be on a similar scale to function correctly.

2. Gradient Descent Convergence: Scaling can help speed up the convergence of gradient descent algorithms by ensuring that all features contribute equally to the gradient.

3. Model Interpretability: It makes it easier to interpret the coefficients in models like linear regression, where the magnitude of the coefficients can be influenced by the scale of the features.

4. Avoiding Dominance: Features with larger scales can dominate the distance calculations in algorithms like k-NN and clustering algorithms, leading to biased results.

### Types of Scaling

There are two primary types of scaling techniques: normalized scaling and standardized scaling.

#### Normalized Scaling

Normalization, also known as Min-Max Scaling, transforms the data to fit within a specific range, usually [0, 1].

Formula:

**Formula:**

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

where:

- $X$ is the original feature value.
- $X_{min}$ is the minimum value of the feature.
- $X_{max}$ is the maximum value of the feature.
- $X'$ is the normalized value.

Pros:

- Keeps the data within a bounded range, making it easier to handle.

Cons:

- Sensitive to outliers, as they can skew the range.

Use Cases:

- When the distribution of the data is not Gaussian (normal).

- When you need to ensure that all features are within a specific range.

#### Standardized Scaling

Standardization transforms the data to have a mean of 0 and a standard deviation of 1. This technique is also known as Z-score normalization.

Formula:

**Formula:**

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

where:

- $X$ is the original feature value.

- $\mu$ is the mean of the feature.

- $\sigma$ is the standard deviation of the feature.

- $X'$ is the standardized value.

Pros:

- Centers the data around 0, making it easier to compare different features.

- Less sensitive to outliers compared to normalization.

Cons:

- Does not bound the data within a specific range.

Use Cases:

- When the data follows a Gaussian (normal) distribution.

- When using algorithms that assume the data is normally distributed, such as linear regression or logistic regression.

### Summary

- Normalization: Transforms data to fit within a specific range (e.g., [0, 1]).

  - Formula:

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

  - Best for: Non-Gaussian distributions, bounded ranges.

  - Sensitivity: Sensitive to outliers.

- Standardization: Centers the data around a mean of 0 and a standard deviation of 1.

  - Formula:

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

  - Best for: Gaussian distributions, algorithms assuming normality.

  - Sensitivity: Less sensitive to outliers compared to normalization.

Both techniques are essential for preparing data for machine learning models, and the choice between them depends on the data distribution and the requirements of the specific algorithm used.

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

VIF (Variance Inflation Factor) can be infinite when there's perfect multicollinearity among predictor variables, meaning one predictor is an exact linear combination of others. This causes the denominator of the VIF formula, 1-R sq 2 to be zero leading to an undefined or infinite VIF

**6.**

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

A Q-Q (quantile-quantile) plot compares the quantiles of a dataset's distribution to a theoretical distribution. In linear regression, it's crucial for assessing the normality assumption of residuals. If the points on the plot fall along the diagonal line, it suggests that the residuals are normally distributed, validating a key regression assumption.