# Linear Regression Subjective Questions

## Assignment-based Subjective Questions

**Ques 1.** From your analysis of the categorical variables from the dataset, what could you infer about their effect on the dependent variable?

**Solution 1:**

Numerous insights can be drawn from the plots which are mentioned below:

- Fall season has the highest rental bike demand.
- Demand for the next year has increased.
- Demand shows continuous monthly growth until June, with September having the highest demand followed by decreasing demand afterward.
- Demand decreases during holidays.
- Weekdays do not provide a clear picture of demand.
- Clear weather situations (weathersit) lead to the highest demand.
- Bike sharing is more in September and less during the end and beginning of the year.

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

**Solution 2:**

Using **drop_first=True** reduces extra columns in dummy variable creation, minimizing correlations among dummy variables. It retains (p-1) dummies to represent p categories. In weathersit, the first column wasn't dropped to preserve information on severe weather situations.

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

**Solution 3:**

Upon examining the pair-plot among the numerical variables, it becomes evident that "temp" and "atemp" exhibit the most substantial correlation (0.63) with the target variable, "cnt."

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

**Solution 4:**

**Residual Analysis:**

- The errors exhibit a normal distribution with a mean of 0, indicating that the model's predictions are centred around the actual values.
- The actual and predicted results follow a similar pattern, suggesting that the model captures the underlying relationships in the data effectively.
- The error terms are independent of each other, which is a crucial assumption for a valid linear regression model.

**R2 Value for Test Predictions:**

The R2 value for predictions on the test data (0.746) is almost identical to the R2 value obtained on the training data (0.758). This high R-squared value indicates that our model performs well even on unseen data, demonstrating its robustness and generalization capability.

**Homoscedasticity:**

The residuals (error terms) exhibit homoscedasticity, with a constant variance across predictions. This implies that the error terms do not show significant variation as the predictor variables change, reinforcing the model's reliability.

**Plot Test vs. Predicted Values:**

The predictions for the test data closely align with the actual values, indicating that the model's performance on unseen data is promising.

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

**Solution 5:**

The top three features are as follows:

- **yr** - This feature shows a positive correlation.

- **season_spring** - This feature shows a negative correlation.
- **weathersit_bad** - This feature also exhibits a negative correlation.

# General Subjective Questions

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

**Solution 1:**

Linear regression is a widely used statistical algorithm for predicting numerical values based on the relationship between one or more independent variables (also known as features, predictors, or input variables) and a single dependent variable (also known as the target variable or output variable). The goal of linear regression is to find the best-fitting linear equation that describes the relationship between the independent variables and the dependent variable, allowing us to make predictions on new data.

The fundamental form of a linear regression model is expressed as:

**$y = a_0 + a_1x_1 + a_2x_2 + ... + a_n*x_n + \varepsilon$**

where:

**y** represents the dependent variable (target variable).

**$a_0$** is the y-intercept, indicating the value of y when all independent variables are 0.

**$a_1, a_2, ..., a_n$** are the coefficients (also known as weights) that represent the slope or impact of each independent variable on the dependent variable.

**$x_1, x_2, ..., x_n$** are the independent variables (features) that influence the target variable.

**$\varepsilon$** represents the random error or residual, which accounts for the variability not explained by the model.

The linear regression algorithm seeks to determine the optimal values of $a_0$, $a_1$, $a_2$,..., $a_n$ that minimize the sum of squared residuals (errors) between the predicted values and the actual values in the training data.

The steps involved in the linear regression algorithm are as follows:

**Data Preparation:** Collect and preprocess the data, ensuring it is clean and well-organized. Divide the data into a training set (used to train the model) and a test set (used to evaluate the model's performance).

**Model Training:** During the training phase, the algorithm adjusts the coefficients a0, a1, a2, ..., an to minimize the error between the predicted and actual values. This process is often achieved through the method of least squares, which finds the line that best fits the data by minimizing the sum of squared residuals.

**Model Evaluation:** Once the model is trained, it is evaluated using the test set. The model's performance is assessed using metrics such as R-squared, Mean Squared Error (MSE), Root Mean Squared Error (RMSE), etc. These metrics help to understand how well the model generalizes to new, unseen data.

**Model Prediction:** After the model is deemed satisfactory, it can be used to make predictions on new data. The input features are fed into the trained model, and it calculates the corresponding dependent variable (target variable).

**Assumptions:** It is important to validate the assumptions of linear regression to ensure the model's reliability. Assumptions include linearity (relationship between variables is linear), homoscedasticity (constant variance of residuals), normality of residuals (residuals follow a normal distribution), and independence of residuals (no autocorrelation).

Linear regression can be extended to handle multiple independent variables (multiple linear regression) and even nonlinear relationships through transformations and feature engineering.

The linear regression algorithm is widely used in various fields, including finance, economics, social sciences, and machine learning, due to its simplicity, interpretability, and effectiveness in modelling relationships between variables. However, it is essential to be cautious about its limitations and suitability for specific datasets, as other more complex algorithms might be necessary in certain cases.


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

**Solution 2:**

Anscombe's quartet is a set of four datasets, each consisting of eleven data points, created by the British statistician Francis Anscombe in 1973. Despite

having vastly different patterns and characteristics, these datasets share nearly identical summary statistics, making them a powerful demonstration of the importance of visualizing data in addition to relying solely on summary statistics.
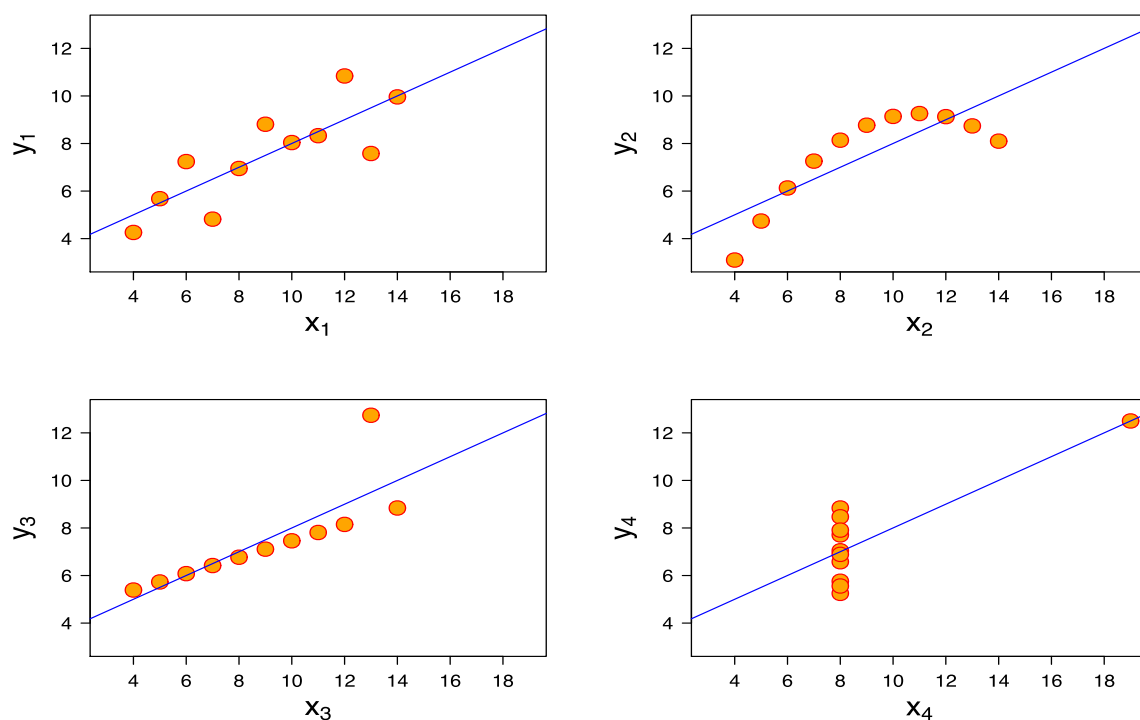
The four datasets within Anscombe's quartet are:

**Dataset I:** This dataset represents a perfect linear relationship. It consists of a series of (x, y) pairs that closely follow a straight line with a slope of approximately 0.5. The summary statistics, including mean, variance, and correlation, are similar to those of the other datasets.

**Dataset II:** This dataset also shows a linear relationship, but it has an outlier at the end that significantly influences the least squares regression line. The outlier causes the regression line to have a steeper slope compared to Dataset I.

**Dataset III:** This dataset appears to have a quadratic relationship. It consists of two distinct groups of (x, y) pairs, where the first group follows a linear relationship, and the second group is quadratic. When looking at summary statistics, it is similar to the previous datasets.

**Dataset IV:** This dataset appears to have no discernible pattern. The (x, y) pairs are randomly scattered with no clear relationship. However, the summary statistics are again similar to those of the other datasets.



All four sets are identical when examined using simple summary statistics, but vary considerably when graphed.

The key takeaway from Anscombe's quartet is that datasets with identical summary statistics can have vastly different patterns and relationships when visually examined. Relying solely on summary statistics, such as mean, variance, and correlation, can be misleading as they do not reveal the true nature of the data. Visualizing the data in graphs and plots is essential for understanding its characteristics and relationships.

Anscombe's quartet serves as a cautionary reminder for data analysts and researchers to always explore and visualize their data before drawing conclusions or making predictions. It highlights the limitations of relying solely on summary statistics and the importance of using data visualization to gain deeper insights and uncover underlying patterns in the data.

**Ques 3.** What is Pearson's R?

**Solution 3:**

Pearson's correlation coefficient, often denoted as "r" or "Pearson's r," is a statistical measure that quantifies the linear relationship between two continuous variables. It was developed by Karl Pearson in the late 19th century and is one of the most widely used measures of correlation in statistics.

Pearson's r ranges from -1 to +1, where:

- **A value of +1 indicates** a perfect positive linear relationship, meaning that as one variable increases, the other variable increases proportionally.
- **A value of -1 indicates** a perfect negative linear relationship, meaning that as one variable increases, the other variable decreases proportionally.
- **A value of 0 indicates** no linear relationship between the two variables.

In other words, Pearson's r measures the strength and direction of the linear relationship between two variables. A positive value indicates a positive correlation, a negative value indicates a negative correlation, and a value close to 0 indicates little to no linear relationship.

To calculate Pearson's correlation coefficient for two variables X and Y with n data points, the formula is:

$$r = \frac{\sum_{i=1}^{n}(X_i - \bar{X})(Y_i - \bar{Y})}{\sqrt{\sum_{i=1}^{n}(X_i - \bar{X})^2}\sqrt{\sum_{i=1}^{n}(Y_i - \bar{Y})^2}}$$

where

r = Correlation Coefficient

Xi = Values of the X-Variable in a sample

Yi = Values of the Y-Variable in a sample

$\overline{X}$ = Mean of the values of the X-Variable

$\overline{Y}$ = Mean of the values of the Y-Variable

Pearson's correlation coefficient is widely used in various fields, including statistics, economics, social sciences, and data analysis, to assess the strength and direction of relationships between variables. It is important to note that Pearson's r only captures linear relationships and may not be appropriate for nonlinear relationships between variables. In such cases, other correlation measures, like Spearman's rank correlation or Kendall's tau, might be more suitable.

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

**Solution 4:**

Scaling, in the context of data preprocessing and feature engineering, refers to the process of transforming the numerical values of features (independent variables) to a specific range or distribution. The goal of scaling is to bring all features to a similar scale or level, which can be beneficial for various reasons in data analysis and machine learning algorithms.

Reasons for performing scaling:

1. **Equalizing Scale:** Features in a dataset may have different scales or units of measurement. Some features might have large value ranges, while others have much smaller ones. Without scaling, features with larger ranges could dominate the model's learning process, leading to biased results.

2. **Improving Convergence:** Many machine learning algorithms, especially those based on gradient descent, converge faster when the features are on a similar scale. Scaling can speed up the optimization process, leading to quicker training times.
3. **Preventing Numerical Instabilities:** Some algorithms, like support vector machines and neural networks, can experience numerical instability or difficulty in calculation if the data is not properly scaled.

There are two common types of scaling techniques:

**Normalized Scaling (Min-Max Scaling):**

- It transforms the features to a fixed range, usually between 0 and 1.
- The formula for Min-Max Scaling is:

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

where:

X = Original value of the feature

Xmax = Maximum value of the feature in the dataset

Xmin = Minimum value of the feature in the dataset

- This method preserves the original distribution and linear relationships between data points but scales them to a common range.

**Standardized Scaling (Z-Score Scaling):**

- It transforms the features to have a mean of 0 and a standard deviation of 1.
- This method standardizes the data, making it suitable for algorithms that assume a Gaussian distribution or those that rely on distance-based metrics.
- The formula for Standardized Scaling is:

$$z = \frac{x - \mu}{\sigma}$$

$$\mu = \text{Mean}$$
$$\sigma = \text{Standard Deviation}$$

## Difference between Normalized Scaling and Standardized Scaling:

- Normalized Scaling brings the features within a fixed range, typically 0 to 1, preserving the original distribution and relative relationships between data points. It is suitable when the data has a bounded range and there are no strong outliers.
- Standardized Scaling, on the other hand, standardizes the data to have a mean of 0 and a standard deviation of 1. It centres the data around 0 and adjusts the scale based on standard deviation. This method is more robust to outliers and is suitable when the data has a Gaussian distribution or when distance-based metrics are used.

Ultimately, the choice of scaling technique depends on the nature of the data and the requirements of the specific machine learning algorithm being used.

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

**Solution 5:**

Yes, it is possible for the Variance Inflation Factor (VIF) to become infinite in certain cases. VIF measures the extent of multicollinearity (high correlation) between predictor variables in a multiple linear regression model. When the VIF is infinite for a particular predictor variable, it indicates that there is a perfect linear relationship between that predictor and other variables in the model.

The formula to calculate the VIF for a predictor variable is:

$$VIF_i = \frac{1}{1 - R_i^2}$$

$R_i^2$ = R-squared value of the linear regression model obtained when the predictor variable in question is regressed against all other predictor variables in the model

Reasons for infinite VIF:

1. **Perfect Collinearity:** Infinite VIF occurs when one or more predictor variables can be perfectly predicted by a linear combination of other predictor variables in the model. In other words, there is an exact linear relationship between the predictor variable and a combination of other predictors. This leads to a situation where the R-squared value is 1 for the regression of that predictor against other predictors, resulting in the denominator of the VIF formula being zero, which leads to an infinite VIF.

2. **Redundant Information:** When a predictor variable is perfectly correlated with a combination of other predictors, it does not add any unique information to the model, and its inclusion can create redundancy.

3. **Dummy Variable Trap:** The dummy variable trap can also lead to infinite VIF when creating dummy variables for categorical predictors. The trap occurs when one dummy variable can be perfectly predicted from the other dummy variables, causing multicollinearity issues.

**Handling infinite VIF:**

Infinite VIF is a serious issue as it indicates severe multicollinearity in the model, which can lead to unreliable regression coefficient estimates and unstable predictions. To address infinite VIF, you need to identify and remove the predictor variable that is causing the multicollinearity problem.

One common approach to handle multicollinearity is to use techniques like stepwise regression or regularization to select a subset of predictors that provide meaningful and unique information while mitigating the impact of

multicollinearity. It is crucial to detect and address multicollinearity in a regression model to ensure the validity and interpretability of the results.

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

**Solution 6:**

A Quantile-Quantile (Q-Q) plot is a graphical tool used to assess whether a given dataset follows a specific theoretical distribution, such as a normal distribution. It compares the quantiles of the observed data with the quantiles that would be expected under the assumed theoretical distribution. The Q-Q plot is especially useful for understanding the distributional properties of a dataset and identifying departures from the assumed distribution.

The steps to create a Q-Q plot are as follows:

- Sort the data in ascending order.
- Calculate the quantiles of the dataset.
- Calculate the expected quantiles based on the assumed theoretical distribution.
- Plot the observed quantiles against the expected quantiles on a scatter plot.
- Use and Importance of Q-Q plot in Linear Regression:

1. **Checking Normality Assumption:** In linear regression, it is often assumed that the error terms (residuals) follow a normal distribution. The Q-Q plot allows us to visually assess whether the residuals have a close-to-normal distribution. If the points in the Q-Q plot approximately follow a straight line, it indicates that the residuals are normally distributed, validating the normality assumption.

2. **Detecting Departures from Normality:** A Q-Q plot can reveal deviations from normality. If the points in the Q-Q plot deviate significantly from a straight line, it suggests that the residuals do not follow a normal distribution. Departures from normality might indicate that the linear

regression model assumptions are violated, and it could lead to biased parameter estimates and inaccurate predictions.

3. **Identifying Outliers:** In addition to assessing normality, a Q-Q plot can help identify outliers in the dataset. Outliers are data points that deviate substantially from the expected quantiles of the theoretical distribution. These outliers might need special consideration in the analysis as they can have a significant impact on the model's performance.

4. **Model Assessment and Improvement:** If the Q-Q plot shows deviations from normality or identifies outliers, it suggests potential issues with the linear regression model. Addressing these issues, such as transforming the data or applying robust regression techniques, can improve the model's fit and predictive capabilities.

5. **Decision-Making in Inference:** When making statistical inferences based on linear regression results, it is crucial to have normally distributed residuals. A Q-Q plot provides visual evidence of whether this assumption holds, and it can guide decisions on the validity and reliability of the statistical inferences.

In summary, a Q-Q plot is a valuable tool in linear regression for assessing the normality assumption of residuals, detecting outliers, and identifying potential issues that can impact the model's accuracy and reliability. It helps researchers and analysts make informed decisions about the adequacy of the model and the validity of the conclusions drawn from the regression analysis.