# Linear Regression Assignment

Name : Anand Bhosale
Batch : C-2

# *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 analysis of the categorical variables provides valuable insights into how they influence the dependent variable, which represents the demand for shared bikes. By examining the coefficients of these variables in the regression model, we can deduce their impact on the dependent variable. When categorical variables have positive coefficients, it indicates a positive effect on bike demand, meaning that higher values of these variables lead to increased bike rentals.

On the other hand, categorical variables with negative coefficients suggest a negative effect, indicating that higher values of these variables are associated with decreased bike rentals. This interpretation enables us to understand the relative importance of different categories within each categorical variable concerning the demand for shared bikes.

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

**Ans :** Using **drop_first=True** during dummy variable creation is important to avoid multicollinearity in the regression model. Multicollinearity occurs when two or more dummy variables are highly correlated, which can lead to unstable coefficient estimates and difficulty in interpreting the model's results. By dropping one of the dummy variables for each categorical feature, we eliminate perfect multicollinearity, ensuring a more robust and reliable model.

**Example: -**

Let's consider a categorical variable "Season" with four categories: Spring, Summer, Fall, and Winter. If we create dummy variables for these categories without dropping the first one, we would have three dummy variables, say "is_Spring," "is_Summer," and "is_Fall." The Winter category is omitted as it is the baseline reference.

Without 'drop_first = True' :

Original Data : Season, Spring, Summer, Fall, Winter

Dummy Variables :

| Is_Spring | Is_Summer | Is_Fall |
|---|---|---|
| 1 | 0 | 0 |
| 0 | 1 | 0 |
| 0 | 0 | 1 |
| 0 | 0 | 0 |

From above, we can see that when all dummy variables are 0, it means the category is "Winter" (the dropped one). This creates perfect multicollinearity between the dummy variables.

Now, With drop_first=True:

Original Data : Season, Spring, Summer, Fall, Winter

Dummy Variables :

| Is_Summer | Is_Fall |
|:---------:|:-------:|
| 0 | 0 |
| 1 | 0 |
| 0 | 1 |
| 0 | 0 |

By dropping the first dummy variable ("is_Spring"), we avoid multicollinearity, and now each dummy variable uniquely represents one category. When both "is_Summer" and "is_Fall" are 0, it means the category is "Winter," and we can interpret the model coefficients more effectively without redundant information.

In summary, using drop_first=True helps in creating a more accurate and interpretable regression model when dealing with categorical variables.

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

**Ans :** Creating a pair plot among the numerical variables is useful for understanding the relationship between each numerical variable and the target variable, which represents the demand for shared bikes. The pair plot allows us to visualize the scatter plots of each numerical variable against the target variable and analyze their trends. The numerical variable that shows the strongest correlation with the target variable is the one where the scatter points closely follow a linear pattern, indicating a substantial positive or negative correlation with the demand for shared bikes. Identifying such influential numerical variables helps in gaining insights into their impact on bike demand.

## 4. How did you validate the assumptions of linear regression after building the model on the training set?

**Ans :** After building a linear regression model on the training set, it is crucial to validate the model's assumptions to ensure its reliability and accuracy. Several common methods are used for assumption validation:

i) Residual Analysis: This involves examining the residuals (the differences between the actual target values and the predicted values) to check if they follow a normal distribution and have constant variance (homoscedasticity). A histogram or a Q-Q plot of the residuals can help assess their normality, while a scatter plot of residuals against predicted values can reveal homoscedasticity.

ii) Multicollinearity Check: To verify if there is multicollinearity among the predictor variables, we can calculate the Variance Inflation Factor (VIF) for each feature. VIF values greater than 5 (sometimes 10) may indicate high multicollinearity.

iii) Linearity Assumption: This assumption suggests that the relationship between the predictor variables and the target variable should be linear. Scatter plots of the numerical predictors against the target variable can help assess this assumption.

iv) Independence of Errors: This assumption assumes that the residuals should be independent of each other. Time-series data may have correlated errors due to autocorrelation, and in such cases, additional techniques like time series cross-validation can be used.

Example: -

In "bike_sharing_dataset" with numerical features like "temperature," "humidity," and "windspeed," and the target variable is "count" (representing the demand for shared bikes). We build a linear regression model on the training set and then validate the assumptions.

1. Residual Analysis: We plot a histogram and a Q-Q plot of the residuals to check for normality. If the residuals are approximately normally distributed, it satisfies the assumption.

2. Multicollinearity Check: We calculate the VIF values for "temperature," "humidity," and "windspeed." If all VIF values are below a certain threshold (e.g., 5), we can assume no significant multicollinearity.

3. Linearity Assumption: We plot scatter plots of "temperature," "humidity," and "windspeed" against "count." If the scatter points follow a relatively linear pattern, the assumption is met.

4. Independence of Errors: If we have time-series data, we might apply time series cross-validation techniques to ensure that the errors are independent.

By validating these assumptions, we can gain confidence in the model's reliability and make more accurate predictions. If any assumption is violated, appropriate measures, such as data transformations or using different modeling techniques, can be applied to address the issue.

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

**Ans :** Based on the final model, the top three features that significantly contribute to explaining the demand for shared bikes can be identified by examining the magnitude of their coefficients. The features with larger absolute coefficients have a stronger impact on the demand for bikes. By comparing the coefficients of all features, we can determine the top three features that have the highest positive or negative coefficients, indicating their significant contribution to explaining the variation in bike demand. These features play a crucial role in influencing the demand for shared bikes and can provide valuable insights for decision-making and business strategy.

# *General Subjective Questions:*

## 1. Explain the linear regression algorithm in detail.

Ans : Linear regression is a widely used statistical algorithm for modeling the relationship between a dependent variable (target) and one or more independent variables (predictors). It aims to find the best-fitting straight line that represents the linear relationship between the independent variables and the dependent variable. This line is referred to as the regression line or the best-fit line.

Detailed explanation of the linear regression algorithm:

**i) Assumptions:** Linear regression assumes that there is a linear relationship between the independent variables and the dependent variable. It also assumes that the errors (residuals) are normally distributed, have constant variance (homoscedasticity), and are independent of each other.

**ii) Simple Linear Regression:** In simple linear regression, there is only one independent variable (predictor) and one dependent variable. The equation of the regression line is represented as:

$y = b0 + b1 * x$

where:
- y is the dependent variable (target).
- x is the independent variable (predictor).
- b0 is the y-intercept (the value of y when x is 0).
- b1 is the slope of the line (the change in y for a unit change in x).

**iii) Multiple Linear Regression:** In multiple linear regression, there are multiple independent variables (predictors) and one dependent variable. The equation of the regression line becomes:

$y = b0 + b1 * x1 + b2 * x2 + ... + bn * xn$

where:
- y is the dependent variable (target).
- x1, x2, ..., xn are the independent variables (predictors).
- b0 is the y-intercept.
- b1, b2, ..., bn are the slopes of the line for each predictor.

**iv) Fitting the Model:** The goal of linear regression is to find the best-fitting regression line that minimizes the sum of the squared differences between the actual target values and the predicted values (residuals). This is usually achieved using the method of least squares.

**v) Coefficient Estimation:** The coefficients b0, b1, b2, ..., bn are estimated during the model fitting process. The process involves finding the values of these coefficients that minimize the sum of squared residuals.

**vi) Model Evaluation:** After fitting the model, it's essential to evaluate its performance. Common evaluation metrics include the R-squared value, which measures the proportion of variance in the target variable explained by the model, and the mean squared error (MSE), which quantifies the average squared difference between actual and predicted values.

**vii) Predictions:** Once the model is trained and evaluated, it can be used to make predictions on new data. By plugging the values of the independent variables into the regression equation, we can obtain the predicted value for the dependent variable.

**viii) Interpretation:** Linear regression allows us to interpret the coefficients of the independent variables. Positive coefficients indicate that an increase in the corresponding predictor variable leads to an increase in the target variable, while negative coefficients suggest a decrease in the target variable.

Overall, linear regression provides a simple yet powerful approach for modeling the relationship between variables and making predictions based on that relationship. However, it is essential to validate the assumptions and assess the model's performance to ensure its reliability and usefulness in practical applications.

## 2. Explain the Anscombe's quartet in detail.

**Ans :** Anscombe's quartet is a set of four datasets in statistics that were created by the British statistician Francis Anscombe in 1973. The unique aspect of these datasets is that they have nearly identical summary statistics (such as mean, variance, correlation, and linear regression coefficients) despite having very different data distributions and patterns. The quartet is often used to highlight the importance of data visualization and to demonstrate that relying solely on summary statistics can be misleading.

Here are the details of the four datasets in Anscombe's quartet:

i) Dataset I:
   - x: 10.0, 8.0, 13.0, 9.0, 11.0, 14.0, 6.0, 4.0, 12.0, 7.0, 5.0
   - y: 8.04, 6.95, 7.58, 8.81, 8.33, 9.96, 7.24, 4.26, 10.84, 4.82, 5.68

   This dataset shows a relatively linear relationship between x and y, with a slight positive correlation. A linear regression line fits the data well.

ii) Dataset II:
   - x: 10.0, 8.0, 13.0, 9.0, 11.0, 14.0, 6.0, 4.0, 12.0, 7.0, 5.0
   - y: 9.14, 8.14, 8.74, 8.77, 9.26, 8.10, 6.13, 3.10, 9.13, 7.26, 4.74

   Similar to Dataset I, this dataset also exhibits a roughly linear relationship between x and y, with a positive correlation. However, it has a different pattern compared to Dataset I.

iii) Dataset III:
   - x: 10.0, 8.0, 13.0, 9.0, 11.0, 14.0, 6.0, 4.0, 12.0, 7.0, 5.0
   - y: 7.46, 6.77, 12.74, 7.11, 7.81, 8.84, 6.08, 5.39, 8.15, 6.42, 5.73

   Dataset III has a clear non-linear relationship between x and y, with an outlier point that significantly affects the linear regression line.

iv) Dataset IV:
   - x: 8.0, 8.0, 8.0, 8.0, 8.0, 8.0, 8.0, 19.0, 8.0, 8.0, 8.0
   - y: 6.58, 5.76, 7.71, 8.84, 8.47, 7.04, 5.25, 12.50, 5.56, 7.91, 6.89

Dataset IV demonstrates how a single outlier can drastically impact summary statistics like the mean and linear regression line.

The purpose of Anscombe's quartet is to emphasize the importance of visualizing data before drawing conclusions. While summary statistics can provide useful insights, they may not reveal the full picture of the underlying relationships and patterns in the data. By plotting the datasets and visually inspecting them, analysts can better understand the data's characteristics and make more informed decisions when performing statistical analyses.


## 3. What is Pearson's R?

**Ans :** Pearson's R, also known as Pearson correlation coefficient, is a measure of the linear correlation between two continuous variables. It quantifies the strength and direction of the linear relationship between variables, ranging from -1 to 1. A value of -1 indicates a perfect negative linear correlation, 0 represents no linear correlation, and 1 signifies a perfect positive linear correlation.
Pearson's R is calculated by dividing the covariance of the variables by the product of their standard deviations. It is commonly used to assess the strength and direction of relationships in various fields, including statistics, social sciences, and finance.


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

**Ans :** Scaling is a data preprocessing technique used in statistics and machine learning to bring all numerical variables or features to a similar scale or range. It involves transforming the values of the variables so that they fall within a specific range, making them comparable and preventing some variables from dominating others due to their larger magnitudes.

Scaling is performed for several reasons:

**Gradient Descent Convergence:** Many machine learning algorithms, especially those based on gradient descent optimization, converge faster when the features are on a similar scale. Scaling prevents the optimization process from being skewed by features with larger values.

**Regularization:** Regularization techniques like L1 and L2 regularization are sensitive to the scale of the features. Scaling helps ensure that the regularization term has a similar impact on all features.

**Distance-based Algorithms:** Algorithms that rely on distance calculations, such as k-nearest neighbors (KNN) or support vector machines (SVM), can be influenced by the scale of features. Scaling makes the distance calculations more meaningful.

**Interpretability:** In some cases, interpreting the model becomes easier when the features are on the same scale, allowing for a better understanding of the feature's impact on the target variable.

**There are two common scaling techniques:**

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

- Normalized scaling scales the features to a fixed range, typically between 0 and 1.
  The formula to perform normalized scaling on a feature x is :

$$x(normalized) = \frac{x - min(x)}{max(x) - min(x)}$$

- This scaling method preserves the original distribution of the data but confines the values within a specified range.

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

- Standardized scaling scales the features to have a mean of 0 and a standard deviation of 1.
  The formula to perform standardized scaling on a feature x is:

$$x(standardized) = \frac{x - mean(x)}{std(x)}$$

- This scaling method transforms the data to have a mean of 0 and a standard deviation of 1, resulting in a standardized distribution.

The key difference between normalized scaling and standardized scaling lies in the range of values the scaled features have. Normalized scaling brings the values within a fixed range (typically 0 to 1), while standardized scaling standardizes the values to have a mean of 0 and a standard deviation of 1. Both scaling techniques are useful, depending on the specific requirements of the machine learning algorithm and the nature of the data being used.

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

**Ans :** The occurrence of infinite VIF (Variance Inflation Factor) values typically indicates perfect multicollinearity among the independent variables. Multicollinearity occurs when two or more independent variables are highly correlated with each other. In such cases, the VIF calculation becomes unstable, resulting in infinite values.

Perfect multicollinearity disrupts the linear regression model because it becomes impossible to estimate the individual effects of highly correlated variables accurately. It leads to numerical instability and inflated standard errors.

To address this issue, multicollinearity should be identified and resolved by removing one of the correlated variables or transforming the variables to reduce correlation.

## 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 (Quantile-Quantile plot) is a graphical tool used to assess whether a dataset follows a particular probability distribution, typically the normal distribution. The plot compares the quantiles of the observed data against the quantiles of the expected distribution.

In linear regression, a Q-Q plot is useful for evaluating the assumption of normality of residuals. By plotting the quantiles of the residuals against the quantiles of the expected normal distribution, we can visually assess if the residuals deviate significantly from normality. If the residuals align closely with the diagonal line in the Q-Q plot, it suggests that the residuals follow a normal distribution, satisfying the assumption of normality in linear regression.

The Q-Q plot helps identify departures from normality, such as skewness or heavy tails, and provides insights into the distributional properties of the residuals, aiding in model evaluation and potential improvements.