

# 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?

## Key Metrics from the Model:

- **R-squared:** 0.800
- **Adjusted R-squared:** 0.797
- **F-statistic:** 286.4

The R-squared value of 0.800 indicates that approximately 80% of the variance in bike rental counts (cnt) is explained by the model. This is a strong indication that the model is a good fit for the data.

## Categorical Variables and Their Effects

### 1. Season (season\_2 and season\_4)

- **season\_2 (summer):**
  - **Coefficient:** 0.0815
  - **t-value:** 7.200
  - **P>|t|:** 0.000
  - **Interpretation:** The coefficient is positive, indicating that, on average, the demand for shared bikes in the summer (season 2) is 0.0815 units higher compared to the reference season (spring, season\_1). The effect is statistically significant given the p-value of 0.000.
- **season\_4 (winter):**
  - **Coefficient:** 0.1265
  - **t-value:** 11.143
  - **P>|t|:** 0.000
  - **Interpretation:** The coefficient is positive, suggesting that, on average, the demand for shared bikes in winter (season 4) is 0.1265 units higher compared to the reference season (spring, season\_1). This effect is highly significant with a p-value of 0.000.
- **Inference:** Both summer and winter have higher bike demand compared to spring, with winter having a more substantial impact.

### 2. Year (yr\_1)

- **Coefficient:** 0.2339
- **t-value:** 25.893
- **P>|t|:** 0.000
- **Interpretation:** The coefficient is positive, indicating that the demand for shared bikes in the year 2019 (yr\_1) is 0.2339 units higher on average compared to the reference year 2018 (yr\_0). This is a strong and significant effect, highlighting the increased popularity of bike-sharing in 2019.
- **Inference:** The bike-sharing demand significantly increased in 2019 compared to 2018.

3. **Month (mnth\_9)**
  - **Coefficient:** 0.0864
  - **t-value:** 5.003
  - **P>|t|:** 0.000
  - **Interpretation:** The positive coefficient indicates that, on average, the demand for shared bikes in September (mnth\_9) is 0.0864 units higher compared to the reference month (January, mnth\_1). This effect is statistically significant.
  - **Inference:** September shows a higher demand for shared bikes compared to January.
4. **Weather Situation (weathersit\_3)**
  - **Coefficient:** -0.2508
  - **t-value:** -9.321
  - **P>|t|:** 0.000
  - **Interpretation:** The negative coefficient suggests that, on average, the demand for shared bikes during bad weather conditions (weathersit\_3, which includes Light Snow, Light Rain + Thunderstorm + Scattered clouds, and Light Rain + Scattered clouds) is 0.2508 units lower compared to the reference weather situation (Clear, Few clouds, Partly cloudy, weathersit\_1). This is a significant effect with a p-value of 0.000.
  - **Inference:** Adverse weather conditions significantly reduce the demand for shared bikes.

## Conclusion

- **Season:** Both summer and winter positively influence bike demand compared to spring, with winter having a more substantial effect.
- **Year:** The year 2019 saw a significant increase in bike demand compared to 2018, highlighting growing popularity.
- **Month:** September has higher bike demand compared to January.
- **Weather Situation:** Bad weather conditions significantly reduce bike demand.

Overall, these categorical variables significantly affect the demand for shared bikes, providing valuable insights for business planning and strategy development for BoomBikes.

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

### 1. Simpler Model Interpretation

When you drop the first category, the model coefficients become easier to interpret. Each dummy variable's coefficient represents the difference in the dependent variable compared to the reference category (the dropped category). This clarity helps in understanding the relative effects of each category.

### 2. Prevention of the Dummy Variable Trap

The dummy variable trap occurs when there are redundant dummy variables, leading to perfect multicollinearity. Dropping the first dummy variable for each categorical feature prevents this trap, ensuring that the model's design matrix is full rank and invertible, which is necessary for accurate parameter estimation.

### 3. Consistency in Statistical Software

Some statistical software or packages might automatically drop the first dummy variable to prevent multicollinearity. By explicitly using `drop_first=True`, you ensure consistency and compatibility across different tools and platforms, avoiding potential errors or warnings.

### 4. Better Performance Metrics

A model with fewer redundant variables can lead to better performance metrics. It can improve the precision of estimated coefficients and enhance the overall predictive performance of the model by focusing only on the necessary predictors.

### 5. Avoiding Multicollinearity

Multicollinearity occurs when one predictor variable can be linearly predicted from others. This issue is relevant in the context of categorical variables represented as dummy variables.

**Example:** Consider a categorical variable `season` with four categories: Spring, Summer, Fall, and Winter

.Without `drop_first=True`, you get four dummy variables:

`season_Spring`, `season_Summer`, `season_Fall`, `season_Winter`

- This leads to perfect multicollinearity because the information in one category is redundant given the others. For instance, if `season_Spring`, `season_Summer`, and `season_Fall` are all 0, it must be Winter.

With `drop_first=True`, you create only three dummy variables:

`season_Summer`, `season_Fall`, `season_Winter`

- The reference category (Spring) is implicitly represented when all the other dummies are 0. This avoids multicollinearity and allows for more stable regression estimates

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

- 'causal' and 'registered' are having highest correlation with target variable 'cnt'
- Other than these 'temp' and 'atemp' are next best

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

- **Normality of Residuals** : verified using normal distribution curve
- **Residual Plot**: Plotted the residuals against the predicted values. residuals are randomly scattered around zero.
- **Linearity** : Plotted the predicted values against the observed values of test data to check for linearity

#### 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 model and the stats that are obtained, top 3 features are

- temp (coef: 0.57) (Positive correlation)
- Weathersit\_3 (coef: -0.2508) (Negative Correlation)
- yr\_1 (coef: 0.23) (Positive correlation)

## General Subjective Questions

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

Linear regression is a fundamental algorithm in machine learning and statistics used for modeling the relationship between a dependent variable (target) and one or more independent variables (predictors). The goal is to find the best-fitting line (or hyperplane in higher dimensions) that predicts the dependent variable based on the values of the independent variables. Here's a detailed explanation of the linear regression algorithm:

### 1. Model Representation

#### Simple Linear Regression

For a single predictor variable  $x$  and a response variable  $y$ , the model can be represented as:  $y = \beta_0 + \beta_1 x + \epsilon$

- $y$  is the dependent variable.
- $x$  is the independent variable.
- $\beta_0$  is the intercept (the value of  $y$  when  $x=0$ ).
- $\beta_1$  is the slope of the line (the change in  $y$  for a one-unit change in  $x$ ).
- $\epsilon$  is the error term (the difference between the observed and predicted values).

## Multiple Linear Regression

For multiple predictor variables  $x_1, x_2, \dots, x_p$ , the model can be represented as:  $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p + \epsilon$ , where  $p$  is the number of predictors.

## 2. Assumptions

Linear regression relies on several key assumptions:

- **Linearity:** The relationship between the predictors and the response is linear.
- **Independence:** Observations are independent of each other.
- **Homoscedasticity:** The residuals (errors) have constant variance at every level of the predictor(s).
- **Normality:** The residuals of the model are normally distributed.
- **No multicollinearity:** The predictors are not highly correlated with each other.

## 3. Objective

The objective is to find the values of the coefficients ( $\beta_0, \beta_1, \dots, \beta_p$ ) that minimize the difference between the observed and predicted values of the dependent variable. This is typically achieved by minimizing the sum of the squared residuals (RSS - Residual Sum of Squares).

## 4. Cost Function

The cost function used in linear regression is the Mean Squared Error (MSE)

## 5. Optimization

To minimize the cost function, the most common method is Ordinary Least Squares (OLS), which calculates the coefficients analytically

## 6. Fitting the Model

1. **Formulate the Design Matrix:** Construct the matrix  $X$  which includes a column of ones for the intercept term and columns for each predictor variable.
2. **Compute the Coefficients:** Use the OLS formula to compute  $\beta$ .
3. **Make Predictions:** Use the computed coefficients to make predictions

## 7. Evaluation

After fitting the model, it is essential to evaluate its performance:

- **R-squared ( $R^2$ ):** Indicates the proportion of the variance in the dependent variable that is predictable from the independent variables..
- **Adjusted R-squared:** Adjusts  $R^2$  for the number of predictors in the model.
- **Mean Absolute Error (MAE):** The average of the absolute errors.
- **Mean Squared Error (MSE):** The average of the squared errors.
- **Root Mean Squared Error (RMSE):** The square root of MSE.

## 8. Assumption Validation

It is crucial to validate the assumptions of linear regression using diagnostic plots:

- **Residual Plot:** Residuals vs. fitted values to check for homoscedasticity.
- **VIF (Variance Inflation Factor):** To check for multicollinearity.

## Conclusion

Linear regression is a powerful and interpretable algorithm for understanding the relationship between variables. However, its effectiveness depends on the validity of its assumptions and the quality of the input data. Proper preprocessing, validation, and diagnostic checks are essential to build a robust linear regression model.

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

Anscombe's quartet is a set of four different datasets that have nearly identical simple descriptive statistics, yet appear very different when graphed. The datasets were constructed by the statistician Francis Anscombe in 1973 to demonstrate the importance of graphing data before analyzing it and to illustrate the effect of outliers and other anomalies on statistical properties.

### Key Features of Anscombe's Quartet

Each dataset in Anscombe's quartet has:

- The same mean of both x and y values.
- The same variance of both x and y values.
- The same correlation between x and y
- The same linear regression line

Despite these identical statistical properties, the datasets differ significantly in structure and distribution, which becomes apparent when they are visualized.

## The Four Datasets

### 1. First Dataset (A)

- Appears as a typical dataset suitable for linear regression.
- The data points are distributed in a roughly linear pattern.
- The linear regression line fits the data well.

### 2. Second Dataset (B)

- The x values are more varied and the y values are more uniform.
- There is a clear nonlinear relationship between x and y.
- A simple linear regression model is not appropriate for this data.

### 3. Third Dataset (C)

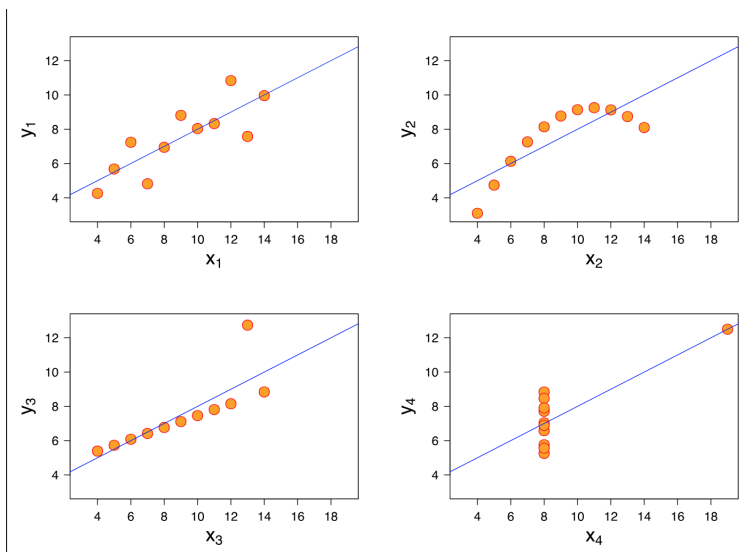
- Most data points have the same x value, except for one outlier.
- The outlier greatly influences the linear regression, skewing the results.
- Without the outlier, the data does not suggest a linear relationship.

### 4. Fourth Dataset (D)

- Most data points are aligned vertically.
- One influential outlier significantly affects the linear regression line.
- The regression line is not representative of the overall data pattern.

## Visual Representation

Here's how the datasets look when plotted:



## Detailed Explanation of Each Dataset

### 1. Dataset A

- **Descriptive Statistics:**
  - Mean of x: 9
  - Mean of y: 7.5
  - Variance of x: 11
  - Variance of y: 4.125
  - Correlation between x and y: 0.816
- **Graphical Insight:**
  - The data points form a linear pattern.
  - The linear regression line fits the data well.

### 2. Dataset B

- **Descriptive Statistics:**
  - Same as Dataset A.
- **Graphical Insight:**
  - The data points form a clear curve, indicating a nonlinear relationship.
  - A linear model is inappropriate.

### 3. Dataset C

- **Descriptive Statistics:**
  - Same as Dataset A.
- **Graphical Insight:**
  - The data points are mostly aligned vertically with one outlier.
  - The outlier skews the linear regression line.

### 4. Dataset D

- **Descriptive Statistics:**
  - Same as Dataset A.
- **Graphical Insight:**
  - The data points are mostly aligned horizontally with one influential outlier.
  - The outlier skews the linear regression line.

## Importance of Anscombe's Quartet

Anscombe's quartet emphasizes several key points:

- **The Importance of Graphical Analysis:** Simple summary statistics can be misleading without graphical analysis. Visualization can reveal patterns, relationships, and anomalies that are not apparent from summary statistics alone.
- **Influence of Outliers:** Outliers can have a significant impact on statistical measures and regression lines. Identifying and understanding outliers is crucial for accurate data analysis.
- **Model Appropriateness:** Different datasets may require different models. A linear model may not be suitable for all datasets, even if summary statistics suggest otherwise.



Anscombe's quartet is a powerful demonstration of why it is essential to visualize data before drawing conclusions from statistical analyses. It shows that datasets with identical statistical properties can have very different structures and relationships, highlighting the need for a comprehensive approach to data analysis that includes both numerical and visual exploration.

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

Pearson's R, also known as the Pearson correlation coefficient or simply the correlation coefficient, is a measure of the strength and direction of the linear relationship between two variables. It is widely used in statistics to quantify the degree of correlation between two continuous variables.

#### Definition and Formula

The Pearson correlation coefficient,  $r$ , is defined as the covariance of the two variables divided by the product of their standard deviations:

$$r = \text{cov}(X, Y) / (\sigma_X \sigma_Y)$$

where:

- $\text{cov}(X, Y)$  is the covariance between variables XXX and YYY.
- $\sigma_X$  is the standard deviation of XXX.
- $\sigma_Y$  is the standard deviation of YYY.

#### Interpretation

- **Value Range:** The value of  $r$  ranges from -1 to +1.
  - **$r=+1$ :** Perfect positive linear relationship. As X increases, Y increases proportionally.
  - **$r=-1$ :** Perfect negative linear relationship. As X increases, Y decreases proportionally.
  - **$r=0$ :** No linear relationship. X and Y are uncorrelated.

#### Significance of Pearson's R

- **Strength of Relationship:** The closer the value of  $r$  is to +1 or -1, the stronger the linear relationship between the variables.
  - **0.9 to 1.0 or -0.9 to -1.0:** Very strong relationship
  - **0.7 to 0.9 or -0.7 to -0.9:** Strong relationship
  - **0.5 to 0.7 or -0.5 to -0.7:** Moderate relationship

- **0.3 to 0.5 or -0.3 to -0.5:** Weak relationship
- **0.0 to 0.3 or 0.0 to -0.3:** Very weak or no relationship
- **Direction of Relationship:**
  - **Positive r:** Indicates that as one variable increases, the other variable also increases.
  - **Negative r:** Indicates that as one variable increases, the other variable decreases.

## Assumptions

- **Linearity:** Pearson's R measures the strength of a linear relationship. It may not be appropriate for non-linear relationships.
- **Homogeneity of Variance:** The variables should have homoscedasticity, meaning the spread of one variable should be roughly the same at all levels of the other variable.
- **Normality:** Ideally, both variables should be normally distributed, especially for smaller sample sizes.

## Applications

Pearson's R is widely used in various fields such as:

- **Social Sciences:** To measure the strength of relationships between variables like income and education level.
- **Medicine:** To understand the correlation between different health metrics.
- **Economics:** To assess the relationship between economic indicators like GDP and unemployment rates.
- **Finance:** To evaluate the relationship between different financial assets.

Pearson's R is a fundamental statistic that provides insight into the linear relationship between two continuous variables. It is a versatile tool that can be applied across multiple disciplines to understand and quantify correlations, guiding further analysis and decision-making. However, it is crucial to remember its limitations and the assumptions underlying its use to ensure accurate interpretation of the results.

## 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 features in a dataset. It is an essential step in preparing data for machine learning models, particularly those that rely on distance measures or gradient descent optimization. Scaling helps ensure that features

contribute equally to the model's learning process and improves the model's performance and convergence speed.

## Why is Scaling Performed?

1. **Equal Contribution:** Different features in a dataset may have different units and ranges. Without scaling, features with larger ranges can dominate the learning process, leading to biased model predictions. Scaling ensures all features contribute equally.
2. **Improved Convergence:** Algorithms that use gradient descent, such as linear regression and neural networks, can converge faster when features are on a similar scale. Unscaled features can lead to slow or unstable convergence.
3. **Distance-Based Algorithms:** Algorithms like k-nearest neighbors (KNN), support vector machines (SVM), and clustering methods (e.g., k-means) rely on distance measures. Scaling ensures that all features are considered equally when computing distances.
4. **Preventing Bias:** Scaling prevents the model from being biased towards features with larger values, ensuring that the importance of all features is assessed correctly.

## Normalized Scaling vs. Standardized Scaling

### Normalized Scaling

Normalization (or Min-Max Scaling) transforms the features to a specific range, typically  $[0, 1]$  or  $[-1, 1]$ .

#### Key Characteristics:

- Scales the data to a fixed range.
- Useful when you want the data to be bounded within a specific range.
- Sensitive to outliers since it depends on the minimum and maximum values.

### Standardized Scaling

Standardization (or Z-score Normalization) transforms the features to have a mean of 0 and a standard deviation of 1.

#### Key Characteristics:

- Centers the data around the mean with a standard deviation of 1.
- Useful when the distribution of the data is approximately normal.
- Less sensitive to outliers compared to normalization.

## Differences Between Normalized and Standardized Scaling

1. **Range:**
  - **Normalization:** Scales the data to a specific range (e.g.,  $[0, 1]$ ).

- **Standardization:** Centers the data around the mean with a unit standard deviation.
- 2. **Outliers:**
  - **Normalization:** Sensitive to outliers, as the min and max values can be affected by extreme values.
  - **Standardization:** Less sensitive to outliers, as it considers the mean and standard deviation.
- 3. **Use Cases:**
  - **Normalization:** Preferred when the data needs to be bounded within a specific range or when the features have different units and you want to bring them to the same scale.
  - **Standardization:** Preferred when the data follows a Gaussian distribution or for algorithms that assume normally distributed data (e.g., linear regression, logistic regression).

Scaling is a crucial preprocessing step in machine learning that ensures all features contribute equally to the model's learning process. Normalized scaling adjusts the data to a specific range, which is useful for distance-based algorithms or when a fixed range is needed. Standardized scaling adjusts the data to have a mean of 0 and a standard deviation of 1, which is useful for algorithms that assume normally distributed data or when the data needs to be centered. Choosing the appropriate scaling method depends on the specific characteristics and requirements of the dataset and the machine learning algorithm being used.

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

The Variance Inflation Factor (VIF) is a measure used to detect the presence and severity of multicollinearity in regression analysis. Multicollinearity occurs when two or more predictor variables in a model are highly correlated, leading to redundancy and instability in the estimated regression coefficients.

### Why the Value of VIF Can Be Infinite

The value of VIF can become infinite under certain conditions, indicating perfect multicollinearity. This happens when one of the predictor variables is an exact linear combination of one or more of the other predictor variables. In other words, the predictor variable can be perfectly predicted by the other predictors in the model.

### Mathematical Explanation

When  $R^2$  becomes 1, VIF becomes infinity

## Causes of Perfect Multicollinearity

1. **Duplicate Columns:** If two columns in the dataset are identical, they will have perfect multicollinearity.
2. **Linear Dependencies:** If one column is a linear combination of other columns (e.g., column A is the sum of columns B and C), this results in perfect multicollinearity.
3. **Dummy Variable Trap:** When creating dummy variables for categorical predictors, if the full set of dummy variables is included in the model, it leads to perfect multicollinearity. To avoid this, one category is usually omitted (using `drop_first=True` in pandas `get_dummies`).

## Implications and Solutions

When VIF is infinite, it indicates a serious problem with the model that must be addressed:

1. **Removing Redundant Predictors:** Identify and remove or combine redundant predictors to eliminate perfect multicollinearity.
2. **Feature Engineering:** Modify the feature set to ensure that no feature is an exact linear combination of others.
3. **Regularization:** Techniques like Ridge Regression can help mitigate the effects of multicollinearity, although they do not solve perfect multicollinearity directly.
4. **Dropping One Dummy Variable:** In the case of the dummy variable trap, always drop one dummy variable to avoid perfect multicollinearity.

## Conclusion

An infinite VIF indicates perfect multicollinearity, where a predictor variable can be perfectly predicted from other predictors. This situation requires immediate attention to adjust the predictor variables and ensure the stability and reliability of the regression model.

## 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 is a graphical tool used to assess whether a dataset follows a particular theoretical distribution, most commonly the normal distribution. It plots the quantiles of the data against the quantiles of the specified theoretical distribution. If the data follows the theoretical distribution, the points on the Q-Q plot will approximately lie on a straight line.

## Construction of a Q-Q Plot

1. **Sort the Data:** Arrange the data in ascending order.

2. **Calculate Quantiles:** Determine the quantiles of the sorted data.
3. **Determine Theoretical Quantiles:** Calculate the corresponding quantiles from the theoretical distribution.
4. **Plot the Points:** Plot the quantiles of the data against the quantiles of the theoretical distribution.

## Interpretation of a Q-Q Plot

- **Straight Line:** If the points fall approximately along a straight line, the data follows the theoretical distribution.
- **Deviations from the Line:** Deviations from the line suggest deviations from the theoretical distribution. For example:
  - **S-shaped Curve:** Indicates that the data has heavier or lighter tails than the theoretical distribution.
  - **Concave/Convex Shape:** Suggests that the data is skewed.

## Use and Importance in Linear Regression

In linear regression, several assumptions must be met for the model to be valid and reliable. A Q-Q plot is particularly useful for assessing the normality of residuals, which is one of these key assumptions.

### Assumptions of Linear Regression

1. **Linearity:** The relationship between the predictors and the response variable is linear.
2. **Independence:** Observations are independent of each other.
3. **Homoscedasticity:** Constant variance of the residuals.
4. **Normality:** Residuals of the model are normally distributed.

## Importance of Q-Q Plot in Linear Regression

1. **Checking Normality of Residuals:** The normality of residuals is crucial for valid hypothesis tests and confidence intervals. If residuals are not normally distributed, the results of statistical tests (e.g., t-tests for coefficients) may be invalid. A Q-Q plot helps visualize whether the residuals follow a normal distribution.
2. **Identifying Outliers:** Q-Q plots can help identify outliers and extreme values that may affect the regression model. Points that deviate significantly from the line in a Q-Q plot indicate potential outliers.
3. **Assessing Model Fit:** By examining the Q-Q plot, one can assess the overall fit of the linear regression model. Deviations from the line indicate areas where the model may not be adequately capturing the data's structure.

## Steps to Use a Q-Q Plot in Linear Regression

1. **Fit the Model:** Fit the linear regression model to the data.

2. **Calculate Residuals:** Compute the residuals (differences between observed and predicted values).
3. **Generate Q-Q Plot:** Plot the quantiles of the residuals against the quantiles of the normal distribution.
4. **Interpret the Plot:** Check if the residuals fall along a straight line to assess normality.