

## **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 :**

From the analysis of categorical variables in the dataset (season, weathersit, holidayand workingday) the following inferences can be drawn regarding their effect on the dependent variable cnt:

### **1. Season**

- **Spring, Summer, Fall, Winter:**
  - Bike demand tends to vary significantly across seasons.
  - **Fall** often shows higher bike rental counts, possibly due to favorable weather conditions.
  - **Winter** may exhibit lower bike rentals due to colder temperatures and harsher weather, reducing outdoor activities.

### **2. Weather Situation (weathersit)**

- **Clear:** Highest bike rentals occur during clear or partly cloudy weather as these conditions are ideal for outdoor biking.
- **Cloudy:** Rentals slightly decline compared to clear weather but remain moderate.
- **Rain:** Rentals decrease significantly due to less favorable conditions.
- **Heavy Rain:** The lowest demand is observed, as such weather greatly discourages biking.

### **3. Holiday**

- **Holiday:** Bike demand may increase or decrease based on the holiday's nature. Demand oriented holidays often correlate with higher bike rentals.
- **Non-Holiday:** Rentals may reflect regular commuting patterns, especially for working days.

### **4. Working Day**

- **Working Day:** Bike rentals are higher during working days, reflecting commuters who use bikes for travel to work.

- **Non-Working Day:** Rentals may drop as fewer people commute. however, demand biking might increase, especially on weekends.

Seasons and weather conditions strongly affect bike demand with favorable conditions driving higher usage. Also Weekdays reflect commuting patterns while weekends/holidays indicate demand trends.

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

**Ans :**

Using `drop_first=True` in dummy variable creation avoids the **dummy variable trap** where including all categories creates redundancy and multicollinearity. By dropping one category, we:

1. **Prevent Multicollinearity:** Ensures predictors are not perfectly correlated.
2. **Set a Baseline:** The dropped category becomes the reference for comparison.
3. **Improve Model Stability:** Makes regression coefficients more interpretable and reliable.

It simplifies the model while maintaining all necessary information.

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

**Ans :**

From the pair-plot among numerical variables (`temp`, `atemp`, `hum`, `windspeed` and `cnt`), **temp** (**temperature**) typically shows the highest positive correlation with the target variable `cnt` (bike rentals). This indicates that as the temperature increases, the demand for bike rentals also rises, likely due to favorable weather conditions for outdoor activities.

### Key Observations:

- **temp vs. cnt:** Strong positive correlation, showing higher bike demand on warmer days.
- **atemp vs. cnt:** Similar trend as `temp`, since it's a "feeling" temperature.

- **hum vs. cnt:** Moderate negative correlation, suggesting high humidity slightly decreases demand.
- **windspeed vs. cnt:** Weak negative correlation, as higher wind speeds discourage biking.

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

**Ans :**

To validate the assumptions of linear regression after building the model:

##### 1. Linearity

- Assumption: The relationship between the independent variables and the dependent variable is linear.
- Validation:
  - Plot the predicted values vs. the actual values or residuals vs. predicted values.
  - Look for a random scatter of residuals. A pattern indicates non-linearity.
  - If non-linearity is detected, consider transformations (e.g, logarithmic or polynomial features).

##### 2. Independence of Errors

- Assumption: Residuals are independent (no autocorrelation).
- Validation:
  - Use the Durbin-Watson test. A value close to 2 indicates no autocorrelation.
  - Particularly important for time-series data.

##### 3. Homoscedasticity

- Assumption: The variance of residuals is constant across all levels of the independent variables.
- Validation:
  - Plot residuals vs. predicted values.
  - Check for a funnel-like shape, which indicates heteroscedasticity (non-constant variance).
  - If heteroscedasticity is present, apply transformations or use weighted regression.

##### 4. Normality of Errors

- Assumption: Residuals are normally distributed.
- Validation:

- Plot a histogram or a Q-Q plot of residuals.
- Perform statistical tests like the Shapiro-Wilk test or Kolmogorov-Smirnov test.
- Non-normality might not critically affect predictions but can impact confidence intervals and hypothesis tests.

## 5. No Multicollinearity

- Assumption: Independent variables are not highly correlated with each other.
- Validation:
  - Calculate the Variance Inflation Factor (VIF) for each independent variable. A VIF > 5 (or 10) suggests multicollinearity.
  - If multicollinearity is detected, consider:
    - Dropping one of the correlated variables.
    - Combining variables (e.g, PCA or feature engineering).

## 6. Model Fit

- Validation:
  - Calculate metrics such as R-squared, Adjusted R-squared, RMSE and MAE on the training and test sets.
  - Check for overfitting or underfitting by comparing training and test performance.

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

**Ans :**

To identify the top 3 features significantly contributing to the demand for shared bikes in the final linear regression model, follow these steps:

### 1. Check Coefficients and P-values

- Examine the **coefficients** of the independent variables in the final model. A larger absolute coefficient value indicates a greater impact on the target variable.
- Use the **p-values** from the regression output to ensure statistical significance. Variables with p-values < 0.05 are considered significant.

### 2. Rank Features by Impact

- After filtering features with significant p-values, rank them by the magnitude of their standardized coefficients (beta coefficients). Standardization is necessary to compare the effect of variables measured on different scales.

### Example Results (Hypothetical):

If the top features identified are:

1. **Temperature (temp)**: Positively correlated, showing demand increases with warmer weather.
2. **Year (yr)**: Indicating a trend of increasing demand from 2018 to 2019.
3. **Working Day (workingday)**: Higher demand on non-holidays or weekdays.

## **General Subjective Questions**

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

Ans :

Linear regression is a fundamental supervised learning algorithm used to establish a relationship between a dependent variable (also known as the target variable) and one or more independent variables (also known as predictors or features). It assumes that this relationship can be expressed as a linear equation. This makes it one of the most widely used methods for regression analysis in statistics and machine learning.

### Key Concepts of Linear Regression

#### 1. Model Equation:

- For a single predictor (x):  $y = \beta_0 + \beta_1 x + \epsilon$
- For multiple predictors ( $x_1, x_2, \dots, x_n$ ):  $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n + \epsilon$

Where:

- y: Dependent variable (target).

- $x_1, x_2, \dots, x_n$  : Independent variables (features).
- $\beta_0$ : Intercept (value of  $y$  when all predictors are 0).
- $\beta_1, \beta_2, \dots, \beta_n$  : Coefficients (weights) representing the contribution of each predictor.
- $\epsilon$ : Error term (captures unexplained variance).

## 2. Linear Assumptions

Linear regression is based on several critical assumptions that must hold for the model to be valid:

- **Linearity**: The relationship between the independent variables (predictors) and the dependent variable (target) is assumed to be linear. This means that any change in the predictors leads to a proportional change in the target variable.
- **Independence of Errors**: The residuals (errors) should be independent of each other, meaning that the error for one observation should not be related to the error for any other observation.
- **Homoscedasticity**: The variance of residuals (errors) should remain constant across all levels of the independent variables. If the spread of residuals changes as a function of the predictors, it is an indication of heteroscedasticity.
- **Normality of Errors**: The residuals are assumed to follow a normal distribution, which is important for conducting hypothesis tests and constructing confidence intervals.
- **No Multicollinearity**: The independent variables should not be highly correlated with each other. Multicollinearity can lead to unstable coefficient estimates and make the model difficult to interpret.

## Steps in the Linear Regression Algorithm

The process of building a linear regression model involves several key steps:

- **Initialize the Model**: Begin by initializing the model with random values for the coefficients ( $\beta_0, \beta_1, \dots, \beta_n$ ).
- **Compute Predictions**: Use the initial values of the coefficients to predict the dependent variable  $y$  for all the data points in the dataset.

- **Calculate the Cost Function:** Measure the model's performance using a cost function such as Mean Squared Error (MSE) or Residual Sum of Squares (RSS). These functions calculate the difference between predicted values and actual values.
- **Optimize the Coefficients:** Adjust the coefficients using optimization algorithms (such as gradient descent) to minimize the cost function and improve model predictions.
- **Evaluate the Model:** Assess the performance of the model using various metrics, including R-squared, Adjusted R-squared, Root Mean Squared Error (RMSE) and Mean Absolute Error (MAE).

### Types of Linear Regression

There are different types of linear regression models based on the number of predictors involved:

- **Simple Linear Regression:** This type of regression involves only a single independent variable ( $x$ ).
- **Multiple Linear Regression:** This type of regression involves multiple independent variables ( $x_1, x_2, \dots, x_n$ ).
- **Regularized Linear Regression:** In this approach, penalties are added to the coefficients to prevent overfitting and improve generalization. Like **Ridge Regression**, **Lasso Regression**

### Strengths of Linear Regression

- **Simplicity:** Linear regression is easy to understand, implement and interpret, making it a great starting point for many types of predictive modeling.
- **Speed:** Linear regression is computationally efficient, especially for small to medium-sized datasets.
- **Baseline Model:** Linear regression often serves as a baseline model for comparison, offering insights into basic relationships in the data before attempting more complex algorithms.

### Limitations of Linear Regression

- **Assumption-Dependent:** Linear regression relies heavily on several assumptions (such as linearity and homoscedasticity), which may not always hold true in real-world data.

- **Sensitive to Outliers:** Outliers or extreme values can disproportionately affect the model's coefficients and predictions, leading to inaccurate results.
- **Limited Expressiveness:** Linear regression can only model linear relationships. To capture non-linear patterns, additional transformations or more complex models (like polynomial regression) may be required.
- **Multicollinearity Issues:** When predictors are highly correlated with each other, it can distort the coefficient estimates and make the model unstable. Multicollinearity can also make it difficult to interpret the influence of individual predictors.

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

**Ans :**

Anscombe's quartet is a set of four data sets that have nearly identical simple descriptive statistics, yet differ significantly in their graphical representations. The purpose of Anscombe's quartet is to demonstrate the importance of visualizing data before performing statistical analysis. Even though these data sets have the same mean, variance, correlation and linear regression fit, the datasets have very different underlying patterns and structures. This illustrates that relying solely on numerical statistics can be misleading and emphasizes the need for graphical analysis.

### Components of Anscombe's Quartet

Anscombe's quartet consists of four data sets, each with 11 pairs of values (x, y). All four data sets have the following identical characteristics:

- **Mean of x:** 9
- **Mean of y:** 7.5
- **Variance of x:** 11
- **Variance of y:** 4.12
- **Correlation between x and y:** 0.82
- **Linear regression line:**  $y = 3 + 0.5x$

Despite these identical statistics, the datasets differ visually, showing that summary statistics alone are insufficient to fully understand the data. Anscombe's quartet is a well-known example in statistics that



highlights the importance of using graphical methods in conjunction with numerical analysis to fully understand data. By visually inspecting data, analysts can detect trends, outliers and non-linear relationships that simple statistics might obscure. The quartet's lesson is clear: "Don't trust statistics without visualizing the data first!"

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

**Ans :**

Pearson's R, also called Pearson's correlation coefficient, is a statistical measure that assesses the strength and direction of the linear relationship between two continuous variables. It indicates how closely one variable changes in relation to another. Pearson's R ranges from -1 to +1.

#### Formula for Pearson's R

The formula to calculate Pearson's correlation coefficient is:

$$r = \frac{n \sum xy - \sum x \sum y}{\sqrt{[n \sum x^2 - (\sum x)^2][n \sum y^2 - (\sum y)^2]}}$$

Where:

- $r$  is the Pearson correlation coefficient.
- $n$  is the number of data points.
- $x$  and  $y$  are the individual data points for the two variables.
- $\sum x$  is the sum of all values in  $x$ .
- $\sum y$  is the sum of all values in  $y$ .
- $\sum xy$  is the sum of the products of corresponding values of  $x$  and  $y$ .
- $\sum x^2$  is the sum of the squares of the values of  $x$ .
- $\sum y^2$  is the sum of the squares of the values of  $y$ .

#### Interpretation of Pearson's R

- **r=1:** A perfect positive linear relationship. As one variable increases, the other also increases in a perfectly linear manner.

- **$r=-1$** : A perfect negative linear relationship. As one variable increases, the other decreases in a perfectly linear manner.
- **$r=0$** : No linear relationship between the variables, though a non-linear relationship may still exist.
- **$0 < r < 1$** : A positive linear relationship, with values closer to 1 indicating a stronger positive correlation.
- **$-1 < r < 0$** : A negative linear relationship, with values closer to -1 indicating a stronger negative correlation.

### Strength of the Correlation

- **0.1 to 0.3** (or -0.1 to -0.3): A weak positive (or negative) correlation.
- **0.3 to 0.5** (or -0.3 to -0.5): A moderate positive (or negative) correlation.
- **0.5 to 1** (or -0.5 to -1): A strong positive (or negative) correlation.

### Considerations

- **Linearity**: Pearson's R measures only linear relationships. It may not accurately reflect the relationship if it is non-linear.
- **Outliers**: Outliers can distort the correlation coefficient, potentially making the relationship appear stronger or weaker than it actually is.
- **Assumptions**: Pearson's R assumes that both variables are normally distributed and have a linear relationship. It is important to validate these assumptions before using this coefficient.

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

Ans :

#### What is Scaling?

Scaling refers to the process of transforming the features of a dataset into a specific range or distribution to make them comparable across different variables. The goal of scaling is to adjust the values of numerical features so they are on the same scale. This is especially important in machine learning

algorithms that are sensitive to the scale of the data, such as those involving distance-based metrics (e.g, K-Nearest Neighbors, Support Vector Machines and Gradient Descent-based algorithms).

### **Why is Scaling Performed?**

Scaling is performed for several key reasons:

1. **To Improve Model Performance:** Many machine learning algorithms perform better when features are on a similar scale. This ensures that no single feature dominates the learning process because of its larger magnitude.
2. **Convergence in Optimization:** Algorithms like gradient descent converge faster when features are scaled, as the learning rate can be uniformly applied to all features.
3. **Equal Weight to All Features:** In models like linear regression, logistic regression, or SVM, features with larger numeric ranges may disproportionately influence the model, causing biased predictions. Scaling puts all features on equal footing.
4. **To Satisfy Assumptions of Certain Algorithms:** Some algorithms (e.g, k-NN, PCA, SVM) assume that the data is scaled properly for accurate distance or similarity measurements.

### **Difference Between Normalized Scaling and Standardized Scaling**

Normalization and standardization are two common methods of scaling, but they differ in how they transform the data:

1. **Normalized Scaling (Min-Max Scaling) :** Normalization is the process of rescaling the features so that they lie within a fixed range, typically  $[0, 1]$  or  $[-1, 1]$ . The data is scaled to a specific range (usually between 0 and 1). It is sensitive to outliers because they can significantly affect the range of the data, thus influencing the normalized values. Often used when the data needs to be bounded or when the distribution is unknown or skewed.
2. **Standardized Scaling (Z-Score Scaling) :** Standardization transforms the data such that it has a mean of 0 and a standard deviation of 1. It centers the data around zero and scales it according to the standard deviation. It does not bound the values within a specific range. The data is centered around zero and the scale is based on the feature's standard deviation. It is not affected by outliers as much as normalization because it uses the standard deviation instead of the range of the data.

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

**Ans :**

The Variance Inflation Factor (VIF) quantifies how much the variance of a regression coefficient is inflated due to multicollinearity, which is the correlation between predictor variables. A high VIF indicates that a predictor is highly correlated with other predictors, leading to unstable and unreliable coefficient estimates.

The VIF can become infinite in the following cases:

Perfect Multicollinearity:

- This happens when one predictor is an exact linear function of another predictor.

Linear Dependence Between Predictors:

- When one predictor is a linear combination of other predictors in the model, the regression model cannot uniquely estimate the coefficients.

Infinite VIF suggests the model is unstable because it cannot distinguish the individual contributions of the correlated predictors. This leads to unreliable coefficient estimates. An infinite VIF is a clear sign of multicollinearity, which needs to be addressed to ensure the model is valid.

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

**Ans :**

A **Quantile-Quantile (Q-Q) plot** is a graphical tool used to assess if a dataset follows a particular theoretical distribution, commonly the **normal distribution**. The plot compares the quantiles of the data distribution against the quantiles of the chosen reference distribution (usually normal).

In a Q-Q plot:

- The **x-axis** represents the quantiles of the theoretical distribution (e.g, normal distribution).
- The **y-axis** represents the quantiles of the observed data.

If the data points lie approximately along a straight line (45-degree line), it suggests that the data follows the theoretical distribution. Deviations from this line indicate discrepancies between the data's distribution and the reference distribution.

### **Use of a Q-Q Plot in Linear Regression**

In linear regression, one of the key assumptions is that the **residuals (errors)** are normally distributed. The Q-Q plot is a useful tool to validate this assumption by visually inspecting whether the residuals of the regression model follow a normal distribution.

#### **Steps to Use Q-Q Plot in Linear Regression:**

1. **Fit a linear regression model:** Perform the regression analysis on the dataset.
2. **Extract the residuals:** Calculate the residuals, which are the differences between the observed values and the predicted values.
3. **Plot the residuals:** Create a Q-Q plot of the residuals against the normal distribution.
4. **Interpret the plot:**
  - If the points closely follow the straight line, it suggests that the residuals are approximately normally distributed, validating the normality assumption.
  - If the points deviate significantly from the straight line, especially in the tails, it suggests that the residuals are not normally distributed, violating the normality assumption.

#### **Importance of a Q-Q Plot in Linear Regression**

1. **Checking Normality of Residuals:**
  - The Q-Q plot helps assess if the residuals follow a normal distribution, which is an important assumption in linear regression. Normally distributed residuals indicate that the model is appropriate for the data and that inference (e.g, confidence intervals, hypothesis tests) is reliable.
2. **Validating Assumptions:**
  - Normality of residuals is essential for the accuracy of statistical tests, such as t-tests and F-tests, which rely on the assumption of normality to draw valid conclusions. The Q-Q plot provides a visual way to check this assumption.
3. **Identifying Outliers:**
  - If the residuals deviate from the normal distribution, particularly in the tails of the plot, it may indicate the presence of outliers or influential data points that are affecting the model. This can guide further investigation into the data.
4. **Model Diagnostics:**

- A Q-Q plot is part of the broader diagnostic process in linear regression. It helps identify potential issues with the model, such as skewed residuals, non-constant variance (heteroscedasticity), or violations of linearity. Addressing these issues improves the robustness and predictive power of the model.

#### **5. Improving Model Assumptions:**

- If the residuals do not follow a normal distribution, transformations on the dependent or independent variables (e.g, log transformation) may help achieve normality, improving the validity of the model.