

Assignment-based Subjective Questions

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

Categorical variables such as season, yr, mnth, holiday, weekday, workingday, and weathersit were analyzed.

Season:

- Rentals are relatively consistent across the four seasons (spring, summer, fall, winter). There is no significant drop or spike, indicating that seasonality may not have a drastic impact on the overall number of rentals, but other seasonal factors could be influencing specific periods.

Weathersit (Weather Situation):

- Most rentals occur under 'Clear' weather conditions, followed by 'Mist'. 'Light Snow or Rain' conditions see a sharp drop in rentals. This suggests that adverse weather conditions discourage bike rentals, with clear weather being the most favourable.

Holiday:

- The majority of rentals happen on non-holidays (0). Very few rentals occur on holidays (1), suggesting that people are less likely to rent bikes on holidays, possibly due to alternative leisure activities or reduced commuting needs.

Workingday:

- More rentals occur on working days (1) compared to non-working days (0). This indicates that a significant portion of bike rentals might be for commuting purposes during workdays.

Weekday:

- Rentals are fairly evenly distributed across all days of the week, with no significant peaks or troughs. This might suggest a steady demand for bike rentals throughout the week, without a strong preference for any particular day.

Month:

- The number of rentals is relatively consistent across different months, with a slight decrease in rentals during the summer months (July and August). This could be due to higher temperatures making biking less comfortable.
-

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

Using drop_first=True helps avoid the dummy variable trap, where the model might suffer from multicollinearity if all categories of a categorical variable are included. By dropping the first category, we prevent this issue and ensure that the model remains stable and interpretable.

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

The pair-plot reveals that temp (temperature) has the highest positive correlation with the target variable cnt (total bike rentals).

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

The assumptions of Linear Regression were validated through several steps:

- **Linearity:** Checked by plotting residuals vs. fitted values to ensure no clear patterns.
 - **Normality:** Residuals were plotted in a histogram or Q-Q plot to check for normal distribution.
 - **Homoscedasticity:** Examined through residuals vs. fitted values to ensure constant variance.
 - **Multicollinearity:** Variance Inflation Factor (VIF) was calculated to ensure no high correlation among predictors.
-

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

The top 3 features based on the final model are:

- temp (temperature) with the highest positive coefficient, indicating it strongly influences bike rentals.
 - yr_2019, indicating a significant increase in rentals in 2019.
 - humidity with a significant negative impact, showing that higher humidity decreases bike rentals.
-

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 denoted as Y) and one or more independent variables (denoted as X). The algorithm tries to fit a line (in the case of simple linear regression) or a hyperplane (in the case of multiple linear regression) that best represents the data. The equation for a simple linear regression model is:

$$Y = b_0 + b_1X + \epsilon$$

Where:

- Y is the dependent variable.
- X is the independent variable.
- b_0 is the intercept (the value of Y when $X=0$).
- b_1 is the slope (the change in Y for a one-unit change in X).
- ϵ (epsilon) is the error term (the difference between the observed and predicted values).

The algorithm finds the best-fitting line by minimizing the sum of the squared differences between the observed values and the values predicted by the line (this method is known as **Ordinary Least Squares** or OLS).

Example: Suppose we have data on house prices (Y) and their sizes (X). A linear regression model could help you predict the price of a house based on its size.

2. Explain the Anscombe's quartet in detail.

Anscombe's quartet consists of four different datasets that have nearly identical simple descriptive statistics (mean, variance, correlation, etc.) but appear very different when graphed. These datasets were constructed by Francis Anscombe in 1973 to demonstrate the importance of data visualization before analysis.

Key Points:

- **Dataset 1:** Represents a simple linear relationship with a small amount of noise.
- **Dataset 2:** Has a clear non-linear relationship, yet the correlation is the same as Dataset 1.
- **Dataset 3:** Contains a single outlier, which heavily influences the correlation.
- **Dataset 4:** All the points lie on a vertical line except for one, drastically skewing the correlation.

Example: Anscombe's quartet shows that relying only on statistical measures can be misleading; visualizing data is crucial to understanding its true structure.

3. What is Pearson's R?

Pearson's R (Pearson's correlation coefficient) is a measure of the linear correlation between two variables X and Y. Its value ranges between -1 and 1:

- **1**: Perfect positive linear correlation.
- **0**: No linear correlation.
- **-1**: Perfect negative linear correlation.

The formula for Pearson's R is:

$$r = \frac{\sum(X_i - \bar{X})(Y_i - \bar{Y})}{\sqrt{[\sum(X_i - \bar{X})^2 * \sum(Y_i - \bar{Y})^2]}}$$

Example: If the Pearson's R between study hours and exam scores is 0.85, it indicates a strong positive relationship, meaning as study hours increase, exam scores tend to increase.

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

Scaling is the process of adjusting the range of features in your dataset to ensure that all features contribute equally to the model. It's important because many machine learning algorithms (like gradient descent-based methods) are sensitive to the scale of data.

- **Normalization (Min-Max Scaling):** Rescales the data to a fixed range, usually [0, 1]. The formula is:

$$X_{scaled} = \frac{X - X_{min}}{(X_{max} - X_{min})}$$

- **Standardization (Z-score Scaling):** Rescales the data to have a mean of 0 and a standard deviation of 1. The formula is:

$$X_{scaled} = \frac{(X - \mu)}{\sigma}$$

5. You might have observed that sometimes the value of VIF is infinite. Why does this happen? (3 marks)

Variance Inflation Factor (VIF) is used to detect multicollinearity in a regression model. A VIF value becomes infinite when there is perfect multicollinearity, meaning one predictor variable is a perfect linear combination of one or more other predictor variables.

Reason: When multicollinearity is perfect, the denominator in the VIF formula ($1 - R^2$) becomes zero, leading to an infinite VIF value.

Example: If two variables, say X1 and X2, are perfectly correlated (e.g., $X1 = 2 * X2$), the VIF for these variables will be infinite.

6. What is a Q-Q plot? Explain the use and importance of a Q-Q plot in linear regression. (3 marks)

A **Q-Q plot** (Quantile-Quantile plot) is a graphical tool to compare the distribution of a dataset to a theoretical distribution (often a normal distribution). It plots the quantiles of the data against the quantiles of the theoretical distribution.

Use in Linear Regression: In linear regression, a Q-Q plot is used to check if the residuals (errors) follow a normal distribution. Normally distributed residuals are an assumption for many inferential statistics, like hypothesis testing.

Example: If residuals from a linear regression model lie on a straight line in the Q-Q plot, it indicates that they are normally distributed. Deviations from this line suggest non-normality, potentially signalling issues like outliers or a need for transformation.