

1)

The summary statistics for `weathersit` likely indicate that bike rentals are affected by weather conditions. For example, rentals are likely higher on clear days compared to days with heavy rain or snow.

The boxplot and summary statistics for `yr` likely show that bike rentals have increased from 2018 to 2019, reflecting growing popularity or expanded service.

2)

Using `drop_first=True` during the creation of dummy variables is important to avoid the scenario that leads to multicollinearity in the dataset.

3)

`atemp` (feeling temperature) has the highest correlation with the target variable

4)

Not sure

5)

The top 3 features contributing towards explaining the demand for shared bikes are

1. `yr_2019`: year 2019, suggesting that bike demand increased in 2019 compared to 2018.
2. `temp`: temperature, with higher temperatures associated with higher bike demand.
3. `atemp`: feeling temperature, which also positively correlates with higher bike demand.

General Questions:-

1)

Linear regression is a fundamental algorithm in machine learning and statistics used to model the relationship between a dependent variable and one or more independent variables. It aims to find the best-fitting linear equation that can predict the dependent variable based on the values of the independent variables.

Types of Linear Regression

1. **Simple Linear Regression**: Involves one dependent variable and one independent variable.
2. **Multiple Linear Regression**: Involves one dependent variable and multiple independent variables.

2)

Anscombe's quartet is a set of four distinct datasets that have nearly identical simple descriptive statistics, yet have very different distributions and appear very different when graphed.

Here are the four datasets:

1. **Dataset I:** A typical linear relationship with some random noise.
2. **Dataset II:** A perfect linear relationship with one outlier.
3. **Dataset III:** A dataset where the linear relationship is only valid due to one influential outlier.
4. **Dataset IV:** A dataset with a vertical relationship caused by a single influential point.

3)

Pearson's R, also known as Pearson's correlation coefficient, is a statistical measure that quantifies the strength and direction of a linear relationship between two continuous variables. It is denoted by r and ranges from -1 to 1.

4)

Scaling transforms data to a standard range or distribution, crucial for machine learning algorithms. It improves model performance, ensures features contribute equally, and speeds up convergence.

- **Normalization (Min-Max Scaling)** scales data to a specific range, which is useful for ensuring all features are within the same range.
- **Standardization (Z-score Normalization)** scales data to have a mean of 0 and standard deviation of 1, which is useful for algorithms that assume normally distributed data or when handling data with outliers.

5)

The Variance Inflation Factor (VIF) can become infinite when there is perfect multicollinearity among the independent variables.

6)

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