

General Subjective Questions:

1. EXPLANATION OF LINEAR REGRESSION ALGORITHM:

Linear regression endeavors to elucidate the relationship between independent and dependent variables using a straight line, applicable solely to numerical variables. The algorithm follows these steps:

- The dataset is partitioned into test and training data.
- Training data is segregated into feature (independent) and target (dependent) datasets.
- A linear model is fitted using the training dataset, often utilizing gradient descent algorithm to determine the coefficients of the best fit line by minimizing a cost function (e.g., residual sum of squares).
- In the case of multiple features, the predicted variable is a hyperplane.
- The predicted variable is then compared with the test data, and assumptions are evaluated.

2. EXPLANATION OF ANSCOMBE'S QUARTET:

Anscombe's quartet comprises four datasets sharing similar simple descriptive statistics but demonstrating distinct distributions when graphically visualized. These datasets challenge the reliance solely on simple statistics for data interpretation, emphasizing the importance of graphical analysis to discern subtle nuances in data patterns.

3. WHAT IS PEARSON'S R:

Pearson's R quantifies the strength of association between two variables, calculated as the covariance divided by the product of their standard deviations. Its values range from +1 to -1, indicating perfect positive linear correlation, no correlation, and perfect negative linear correlation, respectively

4. WHAT IS SCALING? WHY IS SCALING IS PERFORMED? WHAT IS THE DIFFERENCE BETWEEN NORMALIZED SCALING AND STANDARDIZED SCALING?

Scaling standardizes variables to a specific range, a preprocessing step in linear regression. Two common types are:

- Normalized Scaling: Scales variables to a range of 0–1, suitable for non-Gaussian distributions.
- Standardized Scaling: Centers variables around mean with a unit standard deviation, suitable for Gaussian distributions.

5. YOU MIGHT HAVE OBSERVED THAT SOMETIMES THE VALUE OF VIF IS WHY DOES THIS HAPPEN?

VIF (Variance Inflation Factor) becomes infinite when the R square is 1, indicating perfect correlation between features. - $VIF_i = 1/(1 - R_i^2)$

6. WHAT IS A Q-Q PLOT? EXPLAIN THE USE AND IMPORTANCE OF A Q-Q PLOT IN LINEAR REGRESSION?

A Q-Q plot compares two sets of quantiles, serving to check if they originate from the same distribution, aiding in visual data assessment.

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?

Analysis of categorical variables suggests several insights regarding their impact on the dependent variable (Count), including seasonal variations, trends over time, holiday effects,

2. WHY IS IT IMPORTANT TO USE DROP_FIRST=TRUE DURING DUMMY VARIABLE CREATION?

Setting drop_first=True ensures that a variable with n levels can be represented by n-1 dummy variables, simplifying the representation while avoiding multicollinearity issues.

3. LOOKING AT THE PAIR-PLOT AMONG THE NUMERICAL VARIABLES, WHICH ONE HAS THE HIGHEST CORRELATION WITH THE TARGET VARIABLE?

Among the numerical variables, 'temp' exhibits the highest correlation coefficient of 0.63 with the target variable.

4. HOW DID YOU VALIDATE THE ASSUMPTIONS OF LINEAR REGRESSION AFTER BUILDING THE MODEL ON THE TRAINING SET?

Assumptions of linear regression are validated post-model building by examining the distribution of residuals, which should ideally form a normal distribution with a mean of 0.

5. BASED ON THE FINAL MODEL, WHICH ARE THE TOP 3 FEATURES CONTRIBUTING SIGNIFICANTLY TOWARDS EXPLAINING THE DEMAND OF THE SHARED BIKES?

The top three features significantly contributing to bike demand in the final model are identified as 'Holiday', 'Temp', and 'Humidity'.