

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?

A) The fitted regression equation is:

$$\text{cnt} = 0.4028 + 0.5272 \cdot \text{atemp} - 0.1545 \cdot \text{spring} - 0.0974 \cdot \text{holiday} - 0.2195 \cdot \text{snow} - 0.2285 \cdot \text{hum}$$

Based on the analysis and the final fitted equation, the categorical variables have an inverse relationship with demand.

2. Why is it important to use `drop_first=True` during dummy variable creation?
A) It is useful in whether to drop or not drop the first category of categorical variable encoded.
3. Looking at the pair-plot among the numerical variables, which one has the highest correlation with the target variable?
A) The casual and registered user count has very strong relationship with demand count. But, this is obvious. In terms of variables having significant predictive power, 'temp' and 'atemp' have highest correlation with demand at close to 0.63.
4. How did you validate the assumptions of Linear Regression after building the model on the training set?
A) Validated Autocorrelation, and normality based on DW and JB test statistics from `lr.summary()` and multicollinearity through VIF.
5. Based on the final model, which are the top 3 features contributing significantly towards explaining the demand of the shared bikes?
A) Top 3 features: 'atemp', 'snow' ('weathersit') and 'hum'

General Subjective Questions

1. Explain the linear regression algorithm in detail.
A) Linear regression is based on assumption of linear relationship between target and independent variables. The algorithm finds the best fit line to explain the relationship between the dependent and independent variables. The best fit line is identified by minimizing the error (the distance between the points on line and actual values), i.e. Actual – Predicted. These errors are termed as residuals. There are several assumptions related to the residuals for linear regression which includes constant variance of residuals, normal distribution of residuals, absence of serial correlation among residuals.
2. Explain the Anscombe's quartet in detail.
A) The quartet refers to 4 data points which look very similar (identical) based on summary statistics such as mean, median, etc, but very look very different when plotted. This demonstrates importance data visualization.
3. What is Pearson's R?
A) Pearson's R or Pearson's r correlation tells the strength of linear relationship between two variables. It is the ratio of product of variance of the two variables to the covariance of the two variables.
4. What is scaling? Why is scaling performed? What is the difference between normalized scaling and standardized scaling?
A) Scaling is done to increase interpretability of the model coefficients as scaled variables lead to comparable coefficients which are same scale whether standard or normal. Standardized scaling scales the variable using the standard normal distribution scale, i.e, it uses Z-value with 0 as mean 1 as standard deviation of the data.
5. You might have observed that sometimes the value of VIF is infinite. Why does this happen?
A) When Pearson's r is 1, VIF becomes infinity as it is calculated as $1/(1-r^2)$

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

A) Q-Q plot is visual tool to test normality of distribution of a variable as it plots the theoretical (normal) vs actual distribution of the data points.