

Assignment-Based Subjective Questions

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

Ans. The given assignment I first identified Categorical Variable "Season, yr, month, holiday, weekday, workingday, weathersit" against the target variable 'cnt', and EDA visualization was made. As from the Stats this was very clear that weather situation ~~near~~ median where around 50,000 approximately and similar we could see in 'Season' and 'yr'. And the final model building shows a significant growth of R^2 and Adjusted R^2 for 'yr', 'season' etc.

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

Answer! - Yes this is highly advisable to use drop_list = True, this will avoid redundant feature appear while we create correlation and the reference group. This will avoid to create more dummy variable during dummy creation.

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

Ans → Initially I saw numerical variable "registered" having highest correlation with the target variable "cnt". During data preparation we drop 'register' variable due to its multicollinearity the numerical variables 'atemp' got the highest correlation with the target variable 'cnt'.

Q.4. How did you validate the assumptions of linear regression after building the model on the training set?

Ans. Below are the steps performed to carry out the analysis.

1. Test for normal distribution of errors terms (we call it as residuals) by visualizing plot of the error terms.
2. Add and drop the variables based on each model VIF and p-values to avoid multicollinearity

Q.5. Based on the final model, which are the top 3 features contributing significantly towards explaining the demand of the Shared Bikes?

Ans. Below are the top 3 features which significantly explain the demand.

1. atemp - feeling temperature in celsius
2. yr - year (0: 2018, 1: 2019)
3. Winter - It's subcategory of seasons
(4: Winter)

General Subjective Questions

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

Ans It's machine learning algorithm which is based on supervised learning. This performs regression tasks, and models based on independent variables. This is basically used to find out relationship between different variable and its predictive (forecasting) value. In other or simple word we can say this is against dependent variable and independent variables. If we talk about the linear formula it's $y = mx + c$ where c is intercept and m is slope, however for the multivariable this comes with more generic formula.

y (dependent variable)

x (independent variable)

lets explain the linear regression formula as below.

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n$$

where

y is predictive value

β_0 is the constant term

$\beta_1, \beta_2, \beta_3, \dots, \beta_n$ are the model parameters

x_1, x_2, x_3 are the feature values.

Q.2. Explain the Anscombe's Quartet in detail.

Ans. To be frank, this was unknown term for me however I google and try to learn the concept. This is basically a group of 4 data set which looks similar in terms of average value, average value of y . The variance of x and y and the correlation between them is also similar, hence forth linear regression model looks the same however on the scatter plot they tell different story.

Q.3. What is Pearson's R?

Ans. This is basically the strength of the linear regression between the variables. And it varries between -1 and 1. And its feature is if variables goes up and down together then its co-efficient will always be positive and vice versa. Below I tried to explain the theory

$r = -1$ (perfectly linear with negative slope)

$r = 0$ (means no linear association)

$r > 0 < 5$ (means weak association)

$r > 5 < 8$ (means moderate association)

$r > 8$ (means very strong association)

You can find the formula as below

$$r = N \sum_{xy} - (\sum_x)(\sum_y)$$

$$\frac{\sqrt{[N \sum x^2 - (\sum_x)^2][N \sum y^2 - (\sum_y)^2]}}{N}$$

where

N = number of pairs of Scores

\sum_{xy} = Sum of the Product of Pairs Scores

\sum_x = sum of x scores

\sum_y = sum of y scores

\sum_{x^2} = sum of squared x scores

\sum_{y^2} = sum of squared y scores

Q.4, What is scaling? Why is scaling performed?
What is the difference between Normalized
Scaling and Standardized Scaling?

Ans:- This is the technique to standardize the independent features present in the data in a fixed range, and this is performed during the data processing that handle highly varying magnitude.

We do to rescale the variables so that this can be comparable. Imagine if we don't have comparable then coefficient obtained by fitting regression model will

be very large or small compared to the other Co-efficient. And normalized Scale means between 0 to 1.

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

Ans Very simple answer would be that particular variable have severe collinearity. And this variable can be denoted as linear combination of other variables, which mean square of multiple correlation of any predictor variable with the other predictors approaches unity.

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

Ans As per my understanding Q-Q plot means Quantile-Quantile plot, which help us to do assessment on set of data come from some theoretical theoretical distribution such as normal or exponential.

Q-Q plot is Scatterplot created by two set of Quantiles against each other. If both set of Variables come from same set of ~~distribution~~ distribution, then we would see them forming a line, which may be straight line.