

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?

holiday,windspeed,jul,spring,mist+cloud,light rain/snow are negative coefficient if these increase then it has inverse effect in hiring bike.

yr,temp,sep,winter have positive coefficient which means if 1 unit of this variable increases then count of rental bike will also increase.

The equation of our best fitted line is:

$\text{cnt} = 0.253 + (0.253 \times \text{yr}) - (0.098 \times \text{holiday}) + (0.450 \times \text{temp}) - (0.140 \times \text{windspeed}) - (0.073 \times \text{Jul}) + (0.057 \times \text{sep}) - (0.112 \times \text{spring}) + (0.045 \times \text{winter}) - (0.080 \times \text{Mist + Cloudy}) - (0.285 \times \text{lightrain/snow})$

Why is it important to use drop_first=True during dummy variable creation?

When Dummy variables are created for the independent variables categorical once. if not dropped then multicollinearity happens between dummy variable and independent variable. To avoid redundancy in dataset we have to use drop_first.

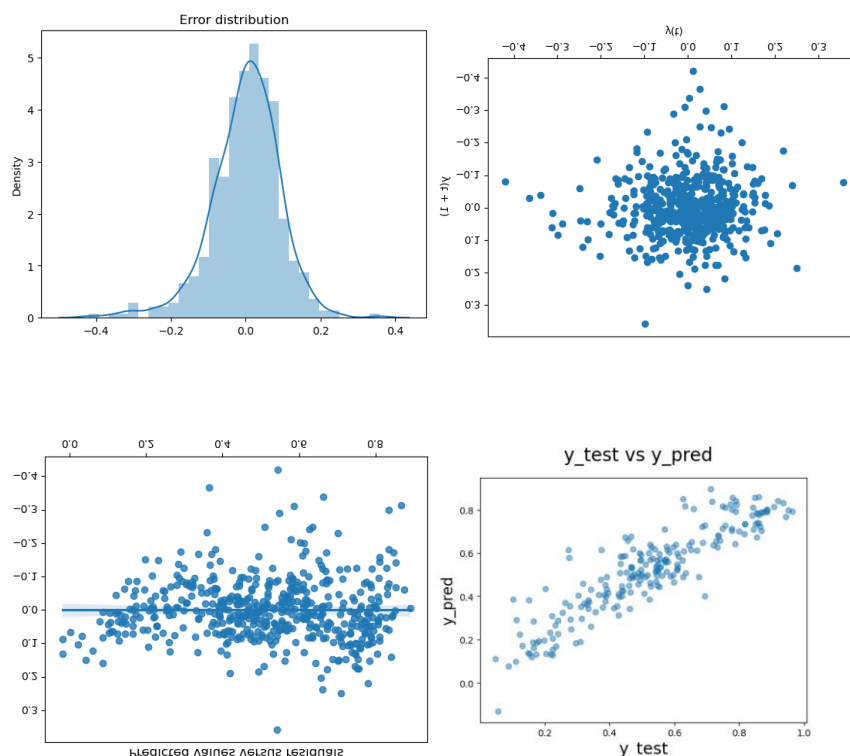
3. Looking at the pair-plot among the numerical variables, which one has the highest correlation with the target variable?

Ans: Temp and atemp

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

Ans: VIF is <5% and $p < 0.05$ which rejects the null hypothesis

Verifying multicollinearity



- Residuals are normally distributed
- Mean Residual errors are independent of each other since the Predicted values vs Residuals plot doesn't show any trend.
- Lagplot of residuals shows no trend. Hence the error terms have constant variance
- Predicted vs observed value plots shows that the model is reasonably accurate.
- Hence, assumptions of Linear Regression are satisfied by this model

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

Ans: Temp, yr (positive coeff) and light rain/snow (negative coefficient) are the three which have highest coefficients among rest in this model

General Subjective Questions

- I. 1. Explain the linear regression algorithm in detail.

Ans: The mathematical equation of Linear regression: $Y = \beta_0 + \beta_1 * x$ for simple Linear regression

Y-Target variable, β_0 is constant, β_1 is independent variable, x is the coefficient.

Positive coeff \rightarrow y increases (coefficient of x) with increase in 1 unit x

Neg coeff \rightarrow y decreases (coefficient of x) with increase in 1 unit x

For multi $Y = \beta_0 + \beta_1 * x + \beta_2 * 2 \dots \beta_i * x_i$

This equation that best describes the correlation of the explanatory variables with the dependent variable. This is achieved by fitting a line to the data using least squares

In Linear Regression, Mean Squared Error (MSE) cost function is used, which helps to figure out the best possible values for a_0 and a_1 , which provides the best fit line for the data points

Steps followed while performing Linear regression:

- I. Draw the scatterplot of the variables. Check for linearity or non-linearity of the data and any deviations. If the pattern is non-linear, perform a transformation. If there are outliers, remove them only if there is a non-statistical reason.

- II. Fit the least-squares regression line to the data and check the assumptions of the model by looking at the residual plot and normal probability plot. If the assumptions of the model are not met, a transformation may be needed.
- III. If necessary, transform the data and re-fit the least-squares regression line using the transformed data.
- IV. If a transformation was performed, then go back to step 1. Otherwise, move to step 5.
- V. Once a “good-fitting” model is obtained, write the equation of the least-squares regression line. Consider the standard errors of the estimates, the estimate of standard deviation, and R-squared.
- VI. Find if the explanatory variable is a significant predictor of the target variable by performing a t-test or F-test.

2. Explain the Anscombe’s quartet in detail.

Anscombe's quartet is a group of datasets (x, y) that have the same mean, standard deviation, and regression line, but which are qualitatively different. It is often used to illustrate the importance of looking at a set of data graphically and not only relying on basic statistic properties.

3.What is Pearson’s R?

Pearson's product moment correlation coefficient (r) is given as a measure of linear association between the two variables: r^2 is the proportion of the total variance (s^2) of Y that can be explained by the linear regression of Y on x.

For example:

- Positive linear relationship: In most cases, universally, the income of a person increases as his/her age increases.
- Negative linear relationship: If the vehicle increases its speed, the time taken to travel decreases, and vice versa.

From the example above, it is evident that the Pearson correlation coefficient, r, tries to find out two things – the strength and the direction of the relationship from the sample sizes.

Strength of Association	Coefficient, r	
	Positive	Negative
Small	.1 to .3	-.1 to -.3
Medium	.3 to .5	-.3 to -.5
Large	.5 to 1.0	-.5 to -1.0

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

It is a step of data Pre-Processing which is applied to independent variables to normalize the data within a particular range. It also helps in speeding up the calculations in an algorithm. Most of the times, collected data set contains features highly varying in magnitudes, units and range. If scaling is not done then algorithm only takes magnitude in account and not units hence incorrect modelling. To solve this issue, we have to do scaling to bring all the variables to the same level of magnitude.

It is important to note that scaling just affects the coefficients and none of the other parameters like t-statistic, F-statistic, p-values, R-squared, etc.

Normalization/Min-Max Scaling:

It brings all of the data in the range of 0 and 1. `sklearn.preprocessing.MinMaxScaler` helps to implement normalization in python.

$$\text{MinMax Scaling: } x = \frac{x - \min(x)}{\max(x) - \min(x)}$$

Standardization Scaling:

Standardization replaces the values by their Z scores. It brings all of the data into a standard normal distribution which has mean (μ) zero and standard deviation one (σ).

$$\text{Standardisation: } x = \frac{x - \text{mean}(x)}{\text{sd}(x)}$$

`sklearn.preprocessing.scale` helps to implement standardization in python.

One disadvantage of normalization over standardization is that it loses some information in the data, especially about outliers.

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

If there is perfect correlation, then $VIF = \text{infinity}$. A large value of VIF indicates that there is a correlation between the variables.

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

Q-Q plots are also known as Quantile-Quantile plots, is a graphical tool to help us assess if a set of data plausibly came from some theoretical distribution such as a Normal, exponential or Uniform distribution. Also, it helps to determine if two data sets come from populations with a common distribution to show if two data sets come from the same distribution. may also be used to determine outliers.

