

Assignment-based Subjective

1. From your analysis of the categorical variables from the dataset, what could you infer about their effect on the dependent variable? (3 marks)

1. Fall season has highest demand for rental bikes.
2. When there is a holiday, demand has decreased.
3. Weekday is not giving a clear picture about demand.
4. Demand is showing continuous growth month on month till June. September month has highest demand. After September, demand is decreasing.
5. During September, bike sharing is more. During the end and beginning of year, it is less.

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

`drop_first=True` is important to use, as it helps in reducing the extra columns that gets created during dummy variable creation. Hence it reduces the correlations created among dummy variables. Dropping the first columns as (p-1) dummies can explain p categories. In `weathersit`, first column was not dropped so as not to lose the info about severe weather situation.

3. Looking at the pair-plot among the numerical variables, which one has the highest correlation with the target variable? (1 mark)

`temp` and `attempt` have the highest correlation (0.63) with the target variable (`cnt`).

4. How did you validate the assumptions of Linear Regression after building the model on the training set? (3 marks)

Residual Analysis: Errors are normally distributed with a mean of 0. Actual and predicted result follow the same pattern. The error terms are independent of each other. R2 value for test predictions: R2 value for predictions on test data (0.815) is almost same as R2 value of train data (0.818). This is a good R-squared value, hence we can see our model is performing good even on unseen data (test data) Homoscedastic: We can observe that variance of the residuals (error terms) is constant across predictions. i.e., error term does not vary much as the value of the predictor variable changes. Plot Test vs Predicted value test: The prediction for test data is very close to actuals.

5. Based on the final model, which are the top 3 features contributing significantly towards explaining the demand of the shared bikes? (2 marks)

The top 3 features are:

1. `yr` (positive correlation)
2. `temp` (positive correlation)
3. `weathersit_bad` (negative correlation)

General Subjective Questions

1. Explain the linear regression algorithm in detail. (4 marks) ?

Linear regression is a statistical regression method used for predictive analysis and shows the relationship between the continuous variables. Linear regression shows the linear relationship between the independent variable (X-axis) and the dependent variable (Y-axis). If there is a single input variable (x), such linear regression is called simple linear regression. And if there is more than one input variable, such linear regression is called multiple linear regression. The linear regression model gives a sloped straight line describing the relationship within the variables. The cost function is used to find the accuracy of the mapping function that maps the input variable to the output variable. This mapping function is also known as the Hypothesis function. In Linear Regression, Mean Squared Error (MSE) cost function is used, which is the average of squared error that occurred between the predicted values and actual values.

2. Explain the Anscombe's quartet in detail. (3 marks) ?

Anscombe's Quartet can be defined as a group of four data sets which are nearly identical in simple descriptive statistics, but there are some peculiarities in the dataset that fools the regression model if built. They have very different distributions and appear differently when plotted on scatter plots. This tells us about the importance of visualising the data before applying various algorithms out there to build models out of them which suggests that the data features must be plotted to see the distribution of the samples that can help you identify the various anomalies present in the data like outliers, diversity of the data, linear separability of the data, etc. Also, the Linear Regression can only be considered a fit for the data with linear relationships and is incapable of handling any other kind of datasets.

3. What is Pearson's R? (3 marks)

Pearson correlation coefficient (PCC), also referred to as Pearson's r, the Pearson product-moment correlation coefficient (PPMCC), or the bivariate correlation, is a measure of linear correlation between two sets of data. It is the covariance of two variables, divided by the product of their standard deviations; thus, it is essentially a normalized measurement of the covariance, such that the result always has a value between -1 and 1. The Pearson's correlation coefficient varies between -1 and +1 where: • $r = 1$ means the data is perfectly linear with a positive slope (i.e., both variables tend to change in the same direction) • $r = -1$ means the data is perfectly linear with a negative slope (i.e., both variables tend to change in different directions) • $r = 0$ means there is no linear association • $r > 0 < 0.5$ means there is a weak association • $r > 0.5 < 0.8$ means there is a moderate association • $r > 0.8$ means there is a strong association
Pearson r Formula Here, • r = correlation coefficient • x_i = values of the x-variable in a sample • \bar{x} = mean of the values of the x-variable • y_i = values of the y-variable in a sample • \bar{y} = mean of the values of the y-variable

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

Scaling 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. When we collect data set it 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 must do scaling to bring all the variables to the same level of magnitude. Normalization/Min-Max Scaling: • It brings all of the data in the range of 0 and 1. 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). • One disadvantage of normalization over standardization is that it loses some information in the data, especially about outliers.

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

If there is perfect correlation, then $VIF = \infty$. This shows a perfect correlation between two independent variables. In the case of perfect correlation, we get $R^2 = 1$, which lead to $1/(1-R^2)$ infinity. To solve this problem, we need to drop one of the variables from the dataset which is causing this perfect multicollinearity. An infinite VIF value indicates that the corresponding variable may be expressed exactly by a linear combination of other variables (which show an infinite VIF as well).

6. What is a Q-Q plot? Explain the use and importance of a Q-Q plot in linear regression. (3 marks)

Quantile-Quantile (Q-Q) plot, 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. This helps in a scenario of linear regression when we have training and test data set received separately and then we can confirm using Q-Q plot that both the data sets are from populations with same distributions. Q-Q Plots (Quantile-Quantile plots) are plots of two quantiles against each other. A quantile is a fraction where certain values fall below that quantile. The purpose of Q Q plots is to find out if two sets of data come from the same distribution. A 45 degree angle is plotted on the Q Q plot; if the two data sets come from a common distribution, the points will fall on that reference line. If the two distributions being compared are similar, the points in the Q-Q plot will approximately lie on the line $y = x$. If the distributions are linearly related, the points in the Q-Q plot will approximately lie on a line, but not necessarily on the line $y = x$. Q-Q plots can also be used as a graphical means of estimating parameters in a location-scale family of distributions.