

## Assignment-based Subjective Questions

**Question 1.** From your analysis of the categorical variables from the dataset, what could you infer about their effect on the dependent variable? (Do not edit)

**Total Marks:** 3 marks (Do not edit)

**Answer:** <Your answer for Question 1 goes below this line> (Do not edit)

In the bike sharing dataset, let's consider the effect of the categorical variable 'weathersit' on the target variable 'cnt'. While performing EDA, I visualized the relationship between the categorical variables and the target variable. It was seen that during the weathersit\_3 (Light Snow, Light Rain + Thunderstorm + Scattered clouds, Light Rain + Scattered), decreases in the bike hires numbers by 0.333164 units have been seen. approximately. Similarly, certain inferences could be made by season\_Spring whereas season\_Winter shows reverse trend. Also, during model building on inclusion of categorical features such as yr, season etc we saw a significant change in the value of R-squared and adjusted R-squared. This implies that the categorical features were helpful in explaining a greater proportion of variances in the data sets

---

**Question 2.** Why is it important to use **drop\_first=True** during dummy variable creation? (Do not edit)

**Total Marks:** 2 marks (Do not edit)

**Answer:** <Your answer for Question 2 goes below this line> (Do not edit)

During dummy value creation (dummy encoding) it is advisable to use drop\_first=True, otherwise we will get a redundant feature i.e. dummy variables might be correlated because the first column becomes a reference group during dummy encoding. For example: suppose we have a categorical feature 'male', 'Female' and 'Cross' as df as shown in fig(a) fig(a) If we not use drop\_first=True we will get all 3 variables as dummy that is as shown in the fig.(b) fig.(b) If we use drop\_first=True we will get all 2 variables as dummy that is sufficient to serve our purpose as shown in the fig(c): If it is Female: Female shows the value as 01. If it is Male: Male shows the value as 10. If it is not Female nor Male its Cross 00. fig(c)

---

**Question 3.** Looking at the pair-plot among the numerical variables, which one has the highest correlation with the target variable? (Do not edit)

**Total Marks:** 1 mark (Do not edit)

**Answer:** <Your answer for Question 3 goes below this line> (Do not edit)

The numerical variable 'registered (0.95)' has the highest correlation with the target variable 'cnt', if we consider all the features. But after data preparation, when we drop registered due to multi collinearity the numerical variable 'atemp (0.63)' has the highest correlation with the target variable 'cnt'.

---

**Question 4.** How did you validate the assumptions of Linear Regression after building the model on the training set? (Do not edit)

**Total Marks:** 3 marks (Do not edit)

**Answer:** <Your answer for Question 4 goes below this line> (Do not edit)

After building the model on the training set, I carried out the following analysis: - Assumptions of Linear Regression: 1. There is a linear relationship between X and Y 2. Error terms are normally distributed with mean zero (not X, Y) 3. Residual Analysis of Training Data proves that the Residuals are normally distributed. 4. Hence our assumption for Linear Regression is valid. 5.

Eliminations and inclusion of independent variables into each model based on VIF and p values to avoid multi collinearity.

---

**Question 5.** Based on the final model, which are the top 3 features contributing significantly towards explaining the demand of the shared bikes? (Do not edit)

**Total Marks:** 2 marks (Do not edit)

**Answer:** <Your answer for Question 5 goes below this line> (Do not edit)

As per our final Model, the top 3 predictor variables that influences the bike booking are: 1. Temperature (temp) - A coefficient value of '0.375922' indicated that a unit increase in temp variable increases the bike hire numbers by 0.375922 units. 2. Weather Situation 3 (weathersit\_3) (Light Snow, Light Rain + Thunderstorm + Scattered clouds, Light Rain + Scattered) - A coefficient value of '-0.333164' indicated that, w.r.t Weathersit\_3, a unit increase in Weathersit\_3 variable decreases the bike hire numbers by 0.333164 units. 3. Year (yr) - A coefficient value of '0.232965' indicated that a unit increase in yr variable increases the bike hires numbers by 0.232965 units. So, it's suggested to consider these variables utmost importance while planning, to achieve maximum Booking.

---

## General Subjective Questions

**Question 6.** Explain the linear regression algorithm in detail. (Do not edit)

**Total Marks:** 4 marks (Do not edit)

**Answer:** Please write your answer below this line. (Do not edit)

<Your answer for Question 6 goes here>

. Linear Regression is a machine learning algorithm based on supervised learning. It performs a regression task. Regression models a target prediction value based on independent variables. It is mostly used for finding out the relationship between variables and forecasting. Different regression models differ based on – the kind of relationship between dependent and independent variables, they are considering and the number of independent variables being used. Linear regression performs the task to predict a dependent variable value (y) based on a given independent variable (x). So, this regression technique finds out a linear relationship between x (input) and y(output). Hence, the name is Linear Regression. The linear regression model can be represented by the following equation:  $Y = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$  where, Y is the predicted value  $\theta_0$  is the constant term.  $\theta_1, \dots, \theta_n$  are the model parameters  $x_1, x_2, \dots, x_n$  are the feature values. The goal of regression analysis is to create a trend line based on the data you have gathered. This then allows you to determine whether other factors apart from the amount of calories consumed affect your weight, such as the number of hours you sleep, work pressure, level of stress, type of exercises you do etc. Before taking into account, we need to look at these factors and attributes and determine whether there is a correlation between them. Linear Regression can then be used to draw a trend line which can then be used to confirm or deny the relationship between attributes. If the test is done over a long time duration, extensive data can be collected and the result can be evaluated more accurately.

---

**Question 7.** Explain the Anscombe's quartet in detail. (Do not edit)

**Total Marks:** 3 marks (Do not edit)

**Answer:** Please write your answer below this line. (Do not edit)

<It is a group of four datasets that appear to be similar when using typical summary statistics, yet tell four different stories when graphed. Each dataset contains of eleven (x, y) pairs as follows:- All the summary statistics for each dataset are identical 1. The average value of x is 9. 2. The average value of y is 7.5. 3. The variance for x is 11 and y is 4.12 4. The correlation between x and y is 0.816 5. The line of best for is  $y = 0.5x + 3$ . But the plots tell a different and unique story for each dataset.

- The first scatter plot (top left) appears to be a simple linear relationship, corresponding to two variables correlated where y could be modelled as gaussian with mean linearly dependent on x.

---

**Question 8.** What is Pearson's R? (Do not edit)

**Total Marks:** 3 marks (Do not edit)

**Answer:** Please write your answer below this line. (Do not edit)

<Your answer for Question 8 goes here>

Pearson's R is a numerical summary of the strength of the linear association between the variables. It varies between -1 and +1. If the variables tend to go up and down together, the correlation coefficient will be positive. If the variables tend to go up and down in opposition with low values of one variable associated with high values of the other, the correlation coefficient will be negative.  $r = 1$  means the data is perfectly linear with a positive slope (i.e., both variables tend to change in the same direction)

By Syed Sha Khalid  $r = -1$  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  $r > 0 < 5$  means there is a weak association  $r > 5 < 8$  means there is a moderate association  $r > 8$  means there is a strong association The Pearson product-moment correlation coefficient (or Pearson correlation coefficient, for short) is a measure of the strength of a linear association between two variables and is denoted by  $r$ . Basically, a Pearson product-moment correlation attempts to draw a line of best fit through the data of two variables, and the Pearson correlation coefficient,  $r$ , indicates how far away all these data points are to this line of best fit (i.e., how well the data points fit this new model/line of best fit). The Pearson correlation coefficient,  $r$ , can take a range of values from +1 to -1. A value of 0 indicates that there is no association between the two variables. A value greater than 0 indicates a positive association; that is, as the value of one variable increases, so does the value of the other variable. A value less than 0 indicates a negative association; that is, as the value of one variable increases, the value of the other variable decreases

---

**Question 9.** What is scaling? Why is scaling performed? What is the difference between normalized scaling and standardized scaling? (Do not edit)

**Total Marks:** 3 marks (Do not edit)

**Answer:** Please write your answer below this line. (Do not edit)

<Your answer for Question 9 goes here>

Scaling is a technique to standardize the independent features present in the data in a fixed range. It is performed during the data pre-processing to handle highly varying magnitudes or values or units. It is extremely important to rescale the variables so that they have a comparable scale. If we don't have comparable scales, then some of the coefficients as obtained by fitting the regression model might be very large or very small as compared to the other coefficients.

Normalized scaling means to scale a variable to have values between 0 and 1, while standardized scaling refers to transform data to have a mean of zero and a standard deviation of 1

---

**Question 10.** You might have observed that sometimes the value of VIF is infinite. Why does this happen? (Do not edit)

**Total Marks:** 3 marks (Do not edit)

**Answer:** Please write your answer below this line. (Do not edit)

<Your answer for Question 10 goes here>

The Variance Inflation Factor (VIF) is a measure of colinearity among predictor variables within a multiple regression. It is calculated by taking the ratio of the variance of all a given model's betas divide by the variance of a single beta if it were fit alone. If there is perfect correlation, then  $VIF = \infty$ . A large value of VIF indicates that there is a correlation between the variables. 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).

---

**Question 11.** What is a Q-Q plot? Explain the use and importance of a Q-Q plot in linear regression. (Do not edit)

**Total Marks:** 3 marks (Do not edit)

**Answer:** Please write your answer below this line. (Do not edit)

<Your answer for Question 11 goes here>

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. Few advantages: a) It can be used with sample sizes also b) Many distributional aspects like shifts in location, shifts in scale, changes in symmetry, and the presence of outliers can all be detected from this plot. It is used to check following scenarios: If two data sets. i. come from populations with a common distribution ii. have common location and scale iii. have similar distributional shapes iv. have similar tail behavior

---