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? (3 marks)

Categorical variables in our dataset played important role on dependent variables, as they are situational and affected the human behaviour to make a decision. As in case of working day which we converted in dummy variable affects the public gathering on certain place, and same phenomenon goes with other categorical variable like "holiday" and "weekday".

```
--------------------------------------------------------------------------------
                       coef     std err        t      P>|t|      [0.025     0.975]
--------------------------------------------------------------------------------
const               1849.3111    204.721      9.033    0.000    1447.093    2251.529
yr                  2001.8857     74.139     27.002    0.000    1856.225    2147.547
season_spring      -1170.1280    144.542     -8.095    0.000   -1454.112    -886.144
mnth_jul            -482.8936    147.516     -3.273    0.001    -772.720    -193.067
season_winter        494.7681    110.193      4.490    0.000     278.271     711.265
mnth_sept            483.0006    135.341      3.569    0.000     217.095     748.907
weekday_sun         -335.8335    103.131     -3.256    0.001    -538.457    -133.210
weathersit_bad     -2305.3236    222.524    -10.360    0.000   -2742.519   -1868.128
weathersit_moderate -665.3845     78.683     -8.456    0.000    -819.974    -510.795
temp                3917.7338    284.697     13.761    0.000    3358.388    4477.080
--------------------------------------------------------------------------------
```

## 2. Why is it important to use drop_first=True during dummy variable creation? (2 mark)

When creating dummy variables, the drop_first=True parameter is used to prevent multicollinearity in the dataset. Multicollinearity occurs when two or more independent variables in a regression model are highly correlated, which can cause issues in the model interpretation and affect the stability of the regression coefficients.

Dummy variables are binary variables created to represent categorical data in a regression model. For a categorical variable with n categories, creating n−1 dummy variables is often sufficient to capture the information about the categories. The dropped category serves as a reference category, and the presence of a 1 in one of the remaining dummy variables indicates that the observation belongs to that category.

By setting drop_first=True, one dummy variable is dropped for each categorical feature, leaving n−1 dummy variables. This helps in avoiding the dummy variable trap, a situation where the dummy variables are perfectly correlated, making it impossible for the regression model to estimate the coefficients accurately.

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

As per our pair-plot "registered" numerical variable has the highest correlation with target variable(cnt).

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

Validating the assumptions of linear regression is a crucial step after building the model on the training set. Here are several common approaches for validating the assumptions. We have used the residual analysis technique to validate our assumptions. We have examined the residuals (the differences between the observed and predicted values) for the day data prediction model created using RFE. Below is the plotted scatterplot of residuals against predicted values to check for linearity, and against each predictor variable to identify potential heteroscedasticity. Utilize residual plots, such as a prediction plot and a Q-Q plot, to assess the normality assumption of residuals.

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

Significant variables to predict the demand for shared bikes

1)holiday

2)temp

3) hum

General Subjective Questions

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

Linear regression is a supervised machine learning algorithm used for predicting a continuous outcome variable (dependent variable) based on one or more predictor variables (independent variables).

The interpretability of linear regression is a notable strength. The model's equation provides clear coefficients that elucidate the impact of each independent variable on the dependent variable, facilitating a deeper understanding of the underlying dynamics. Its simplicity is a virtue, as linear regression is transparent, easy to implement, and serves as a foundational concept for more complex algorithms.

Linear regression is not merely a predictive tool; it forms the basis for various advanced models.

Techniques like regularization and support vector machines draw inspiration from linear regression,

expanding its utility. Additionally, linear regression is a cornerstone in assumption testing, enabling

researchers to validate key assumptions about the data.

There are two main types of linear regression:

1) Simple Linear Regression

This is the simplest form of linear regression, and it involves only one independent variable and one

dependent variable. The equation for simple linear regression is:

where:

Y is the dependent variable

X is the independent variable

$\beta_0$ is the intercept

$\beta_1$ is the slope

2) Multiple Linear Regression

This involves more than one independent variable and one dependent variable. The equation for

multiple linear regression is:

where:

Y is the dependent variable

$X_1, X_2, ..., X_p$ are the independent variables

$\beta_0$ is the intercept

$\beta_1, \beta_2, ..., \beta_n$ are the slopes

In summary, linear regression involves representing a relationship between variables using a linear

equation, defining a cost function to measure the model's performance, and optimizing the model

parameters using an iterative process like gradient descent to minimize the cost and obtain the bestfitting

line.

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

Anscombe's quartet is a set of four datasets that have nearly identical simple descriptive statistics but differ significantly when graphed.

It illustrates the importance of visualizing data and not relying solely on summary statistics.

Dataset composition:

• Anscombe's quartet consists of four datasets, each containing 11 data points. Each dataset has two variables: x (independent variable) and y (dependent variable).

Descriptive Statistics:

• Despite having identical or very similar summary statistics (mean, variance, correlation, and linear regression parameters), the datasets exhibit distinct patterns when graphed.

Illustration of the Importance of Visualization:

• Anscombe's quartet highlights the limitation of relying solely on summary statistics. Even if two datasets have similar mean, variance, and other summary measures, their underlying structures may differ.

• By graphing the data, it becomes evident that the datasets have different distributions, relationships between variables, and patterns of variability.

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

Pearson's r, is a statistical measure that quantifies the strength and direction of a linear relationship between two continuous variables. It was developed by Karl Pearson and is widely used in statistics to assess the linear association between two variables.

The Pearson correlation coefficient can take values between -1 and +1:

r =1: Perfect positive linear relationship

r =-1: Perfect negative linear relationship

r =0: No linear relationship

Pearson's correlation coefficient is particularly useful for assessing the linear relationship between variables, but it assumes that the relationship is linear and that the data is approximately normally distributed. If the relationship is not linear, Pearson's "r" may not accurately capture the

association between variables. In such cases, other correlation measures or non-linear regression techniques may be more appropriate.

### 4. What is scaling? Why is scaling performed? What is the difference between normalized scaling and standardized scaling? (3 marks)

Scaling in linear regression refers to the process of normalizing or standardizing the input features of the model. The purpose is to bring all the features to a similar scale, typically by transforming them in a way that their values have similar ranges. This is important because linear regression models are sensitive to the scale of the input features, and features on different scales can impact the performance and convergence of the model.

Scaling is performed for below reasons-

Improves Convergence: Scaling helps the optimization algorithm converge faster. When features are on a similar scale, the optimization process is more efficient, and it can find the optimal coefficients for the features more quickly.

Equalizes Variable Influence: Scaling ensures that all variables contribute to the model fitting process more uniformly. Without scaling, features with larger magnitudes can dominate the learning process, leading to an unbalanced influence on the model.

Facilitates Interpretability: Scaling doesn't affect the interpretation of the coefficients in terms of the feature importance. It just helps the optimization process. The relationships and significance of coefficients remain the same.

Normalized Scaling (Min-Max Scaling):

Scales the features to a specific range, usually between 0 and 1. Normalized scaling is sensitive to outliers because it depends on the range of the data.

Standardized Scaling (Z-score Scaling):

Standardizes the features to have a mean of 0 and a standard deviation of 1. Standardized scaling is less sensitive to outliers because it uses the mean and standard deviation, which are less affected by extreme values.

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

The Variance Inflation Factor (VIF) is a measure used to assess the severity of multicollinearity in a multiple regression analysis. High VIF values indicate that the variance of the estimated regression coefficients is inflated due to collinearity among the predictor variables. However, in case VIF is infinity that is mean the R-Square value is 1, as the VIF has $(1-r^2)$ in Denominator.

In case r^2 is 1 which represent that we are able to predict all the values 100% which is not the ideal case and it is overfitting the model, basically model has memorized all the value which is not the good model.

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

A Q-Q (quantile-quantile) plot is a graphical tool used to assess whether a dataset follows a

theoretical distribution, such as the normal distribution. In the context of linear regression, Q-Q plots are often employed to check the normality assumption of the residuals.

Here's an explanation of the use and importance of Q-Q plots in linear regression:

Use of Q-Q Plot in Linear Regression:

Assumption Checking:

One of the key assumptions in linear regression is the normality of the residuals (the differences between the observed and predicted values). Q-Q plots help visualize whether the residuals follow a normal distribution.

Comparing Distributions:

The Q-Q plot compares the quantiles of the observed residuals with the quantiles expected from a theoretical normal distribution. If the points on the Q-Q plot closely follow a straight line, it indicates that the residuals are approximately normally distributed.