

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)

Interpretation on “Weather Situation” and “Season” Categorical columns:

When the weather is Clear, then the overall (“cnt”) number of rides are high compared with “Mist” and “Light snow / Light Rain” weather.

However, the difference between the number of rides when the weather is Clear vs Mist are very less. But there is a significant drop in number of rides when the weather is “Light snow / Light Rain”

During all the three types of weather conditions like Clear / Mist / Light Snow or Rain:

- It’s clearly evident that Spring season has a lesser number of rides compared with the other three seasons
- And the “Fall” season has highest number of rides compared with the other three seasons
- And the number of rides is increased year by year across all the weather conditions

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

Multicollinearity is when two or more independent variables in a regression model are highly correlated with each other. It may lead to problems in interpreting the coefficients of the model and can make model unstable. When there is categorical variable present with n categories, then it leads to creation of n dummy variables each represents on categorical value. But it has redundant information. Setting drop_first=True ensures we only create n-1 dummy variables. This avoids redundancy. It means, drop_first=True helps prevent multicollinearity.

For Example, *MaritalStatus* categorical column may have two values ‘Married’ or ‘Unmarried’. If we convert into dummy variables, it leads two columns one for *Married* and other for *Unmarried*, both are highly correlated. If Married column indicates 1 then obvious Unmarried column indicates 0. Both columns represent the same meaning causes redundancy. So it is obvious and more meaningful to drop one column for better stable model.

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

“temp” and “atemp” variables have the highest correlation. Then “registered” and “cnt” (target variable) variables are having the second high correlation from the pair plot diagram

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

Here are the main assumptions of Linear Regression:

Linear Relationship: If we plot Scatter Plot on the independent variable vs. the dependent variable and the output results into a rough straight line, it suggests a linear relationship. The linear relationship is observed between predicted values and actual values.

Homoscedasticity (Equal Variance of Residuals): If we scatter plot on Residuals vs. Fitted Values, the residuals on the y-axis and the fitted values (predicted values) on the x-axis and if the spread of the residuals is constant across the range of fitted values, it suggests homoscedasticity. This is overed with scatter plot after on the predicted values vs residuals.

No Multicollinearity: Calculate the correlation matrix between all the independent variables. Look for correlations close to 1 or -1, which indicate strong collinearity. This can be a starting point, but there are more sophisticated methods for detecting multicollinearity.

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 variables with Higher magnitude Coefficients, Lowest P-values will contribute significantly to explain the demand of the shared bokes. From the model parameters, we could see that Casual drives, Working Day (indicates business day), and the year (year by year increase) are contributing **positive significantly** towards the demand of the shared bikes

Similarly, Weather situation Light Snow or Light Rain, Spring season and windspeed are in the order of **negatively** impacting the demand of the shared bikes

Coefficient

	coef	std err	t	P> t	[0.025	0.975]
const	0.1831	0.022	8.286	0.000	0.140	0.226
yr	0.1815	0.007	25.765	0.000	0.168	0.195
mnth	0.0343	0.018	1.924	0.055	-0.001	0.069
weekday	0.0210	0.010	2.134	0.033	0.002	0.040
workingday	0.1864	0.009	20.249	0.000	0.168	0.205
windspeed	-0.0701	0.021	-3.414	0.001	-0.110	-0.030
casual	0.6689	0.025	26.442	0.000	0.619	0.719
spring	-0.1322	0.015	-9.060	0.000	-0.161	-0.104
summer	-0.0364	0.011	-3.311	0.001	-0.058	-0.015
LightSnowRain	-0.1862	0.020	-9.197	0.000	-0.226	-0.146
Mist	-0.0491	0.007	-6.855	0.000	-0.063	-0.035

General Subjective Questions

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

Linear regression is an algorithm in machine learning used for predicting continuous values based on input features. It assumes a linear relationship between these features and the target. Linear regression aims to find a straight line that best fits this data. This line represents the equation of your linear model, which can be used to predict the value of Y for new unseen X values.

There are two main categories of linear regression depending on the number of independent variables:

Simple Linear Regression: This involves only one independent variable (X) to predict the dependent variable (Y).

Multiple Linear Regression: This involves multiple independent variables (X1, X2, X3, etc.) to predict the dependent variable (Y).

The key aspect of linear regression is finding the equation for the best fit line. This is achieved by minimizing the residuals. Residuals represent the vertical distances between each data point and the corresponding point on the fitted line.

There are different techniques to achieve this, but a common approach is the method of least squares. It involves minimizing the sum of squared residuals, essentially finding the line that makes the sum of these squared distances as small as possible.

The equation of a straight line is generally represented in mathematics as: $Y = b_0 + b_1 * X$

where, b_0 is the y-intercept (the point where the line crosses the Y-axis).

b_1 is the slope of the line (represents the change in Y for a unit change in X).

The linear regression algorithm determines the optimal values for b_0 and b_1 that minimize the residuals.

The main assumption is that the relationship between X and Y is linear. If the data exhibits a non-linear pattern, the predictions won't be accurate.

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

Anscombe's quartet is a modal example in statistics created by Francis Anscombe in 1973. It consists of four sets of data, each containing 11 (x, y) data points. These datasets have identical basic summary statistics like mean, variance, standard deviation, correlation coefficient, etc.

But when we plot these data sets, they reveal themselves to be very different visually. This explains the crucial role of data visualization in understanding the underlying patterns in the data. Eventually, it explains the necessity of performing Exploratory Data Analysis (EDA) on the data sets before we build any ML model.

Anscombe's quartet demonstrates that any numerical statistics can be misleading and fail to capture the true essence of the data if we do not visualize the data. By visualizing the data, we can uncover important aspects like:

Non-linear relationships: Even with similar means and correlations, the data points might not follow a straight line, indicating a non-linear relationship between the variables.

Outliers: Summary statistics might not reveal the presence of outliers that can significantly influence the fitted line.

Underlying distributions: The visual representation can show us the spread and distribution of the data points, which can be crucial for choosing appropriate statistical methods.

The Anscombe's quartet explains that if there are four different Data Sets and let us say the respective scatter plots drawn on these Datasets may represent as follows:

Data Set 1: This could be a tight cluster of points following a clear linear trend.

Data Set 2: This might show a curved pattern despite having the same statistics as Set 1.

Data Set 3: This could be a random scatter with one extreme outlier heavily influencing the fitted line (although the statistics remain similar).

Data Set 4: This might depict a straight line, but with all the data points clustered on one side, making the line irrelevant for most predictions.

However, the mean, variance, standard deviation of the above all four Datasets might be same / nearly matching though the data points represent different scatter plots. This concludes that it is essential to perform EDA / visualize the data before building any complex models based on statistics.

Visualize your data: Create scatter plots, histograms, and other visualizations to understand the distribution of your variables and identify potential issues.

Look for outliers: Investigate outliers and assess their impact on the analysis.

Consider non-linear relationships: Don't blindly assume a linear relationship. Explore if transformations or different models might be more suitable.

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

Pearson's R, also known as the Pearson correlation coefficient, is a statistical measure used to quantify the linear relationship between two continuous variables. It represents the strength and direction of that association.

Strength: The value of R ranges from -1 to +1.

Positive R (between 0 and +1): Indicates a positive correlation. As the value of one variable increases, the other variable tends to increase as well. (Stronger positive correlation as the value gets closer to 1).

Negative R (between -1 and 0): Indicates a negative correlation. As the value of one variable increases, the other variable tends to decrease. (Stronger negative correlation as the value gets closer to -1).

Direction: The positive or negative sign indicates the direction of the relationship.

No Correlation: A value of 0 indicates no linear correlation. The changes in one variable are unrelated to the changes in the other.

Pearson's R only measures linear relationships. It won't capture non-linear patterns like curves or exponential trends.

It assumes normally distributed data for both variables. If the data is not normally distributed, the R value might be misleading.

Pearson's R is widely used in various fields to understand relationships between variables. Examples like Analyzing the correlation between stock prices and market movements, Investigating the correlation between environmental factors and plant growth, etc

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

Scaling is a data preprocessing technique where we adjust the values of features within a dataset to a specific range proportionally matching with original data values. This is done to ensure all features contribute equally to the model's training process and improve the overall performance.

Algorithms like gradient descent rely on calculating the difference between predicted and actual values. Features with vastly different scales can cause the updates to become very small for features with lower scales, hindering the convergence process. Scaling helps achieve a more balanced optimization process.

Features with a much larger range of values can dominate features with smaller ranges during model training. Scaling ensures each feature has a similar influence on the model.

Computational power and usage of resources on scaling data is more optimum and efficient.

There are two most common scaling techniques:

Normalization: Normalization typically scales features to a range between 0 and 1 (or -1 and 1). Here's a common approach for min-max normalization:

$$\text{New_Value} = (\text{Old_Value} - \text{Min_Value}) / (\text{Max_Value} - \text{Min_Value})$$

where,

Min_Value is the minimum value of the feature in the dataset.

Max_Value is the maximum value of the feature in the dataset.

This ensures all features are restricted to a specific range. It's useful when you want to bound the data within a known range of 0 and 1.

Standardization: Standardization transforms features to have a zero mean and unit standard deviation. Here's the formula:

$$\text{New_Value} = (\text{Old_Value} - \text{Mean}) / \text{Standard_Deviation}$$

This process centers the data around zero and scales it based on the spread of the data points. It's particularly useful when the features have different units of measurement, and you want the model to focus on the relative differences between data points rather than the absolute values.

Differences:

Normalization: Creates a uniform distribution of values between a specific range (often 0-1 or -1 to 1).

Standardization: Creates a normal distribution with a mean of 0 and a standard deviation of 1.

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

An infinite VIF (Variance Inflation Factor) in linear regression occurs when there's perfect multicollinearity among the independent variables. This means one independent variable can be expressed entirely as a linear combination of the other independent variables.

Reasons for infinite VIF is explained as follows:

VIF Calculation: VIF is calculated as $1 / (1 - R^2)$, where R^2 is the coefficient of determination between a specific independent variable and all the other independent variables combined.

Perfect Multicollinearity: In this case, R^2 between the variable and the others becomes exactly 1. This signifies a perfect linear relationship, where one variable can be perfectly predicted by the others.

Division by Zero: Plugging R^2 of 1 into the VIF formula results in a denominator of zero ($1 - 1 = 0$). Dividing by zero is mathematically undefined, hence the infinite VIF.

Essentially, an infinite VIF indicates that the variable you're looking at is redundant because it carries no unique information. It can be completely recreated using the other independent variables.

Unreliable Coefficients: The regression coefficients for the variables with high VIF (including the one with infinite VIF) become unreliable and difficult to interpret.

Inaccurate Predictions: The model might not be able to make accurate predictions due to the inflated variances of the coefficients.

We can handle High VIFs / Infinite VIFs in the following approach:

Identify Redundant Variables: Analyze the correlation matrix with high correlations or calculate VIFs of independent variables (indicating potential multicollinearity).

Remove Redundant Variables: Remove one or more of the highly correlated variables based on domain knowledge or feature importance analysis. Probably, $VIF > 5$ might be removed, where 5 can be adjusted to a different value based on problem statement.

By addressing high / infinite VIF, we can ensure your regression model provides reliable results.

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 plot, short for quantile-quantile plot, is a graphical tool used to assess the normality of residuals in linear regression. It helps visualize how closely the distribution of the residuals aligns with a theoretical normal distribution.

The steps involved are

- 1) Calculating Quantiles: Both the observed residuals from a linear regression model and the theoretical quantiles from a normal distribution are calculated. Quantiles represent specific percentiles of the data distribution.
- 2) Plotting the Quantiles: The quantiles of the residuals are plotted against the corresponding quantiles of the normal distribution. Ideally, the points should fall roughly along a straight diagonal line.

Interpreting the Q-Q Plot:

Straight Diagonal Line: If the points form a reasonably straight diagonal line, it suggests that the residuals are normally distributed. Deviations from the straight line indicate departures from normality.

Points above the line indicates that the tails of the distribution are heavier than a normal distribution (fatter tails), possibly indicating outliers or skewed data.

Points below the line indicates that the lighter tails in the distribution compared to normal, potentially implying a presence of more data points near the center.

A Q-Q plot provides a visual way to assess normality of residuals and identify potential issues. And Q-Q plot is helpful in identifying Outliers. Deviations from the line in the Q-Q plot can highlight the presence of outliers that might be affecting the model. The Q-Q plot is a valuable tool for diagnosing potential problems with the linear regression model.

By incorporating Q-Q plots into your linear regression analysis, you gain a better understanding of the underlying distribution of residuals.