

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)

Equation is

$$\text{cnt} = 4491.30 + 1069.78 \times \text{yr} + 370.72 \times \text{september} - 257.79 \times \text{november} - 263.97 \times \text{may} + 318.96 \times \text{june} - 344.55 \times \text{Cloudy} - 367.86 \times \text{july} - 216.54 \times \text{february} - 265.03 \times \text{december} + 260.30 \times \text{july} + 357.72 \times \text{winter} + 365.59 \times \text{august} - 454.45 \times \text{light snow/rain} - 391.91 \times \text{spring}$$

The categorical variables (months, seasons, and weather) have a significant impact on cnt. The model highlights the importance of understanding seasonal trends, weather effects, and yearly growth to optimize bike rental operations and strategies.

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)

When creating dummy variables (also known as one-hot encoding) for categorical variables, it's important to use the `drop_first=True` parameter in certain cases to avoid multicollinearity and to ensure our regression model works correctly.

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)

Year(yr)

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 a Linear Regression model, it's essential to validate the assumptions to ensure the model's reliability and accuracy. The key assumptions of Linear Regression are:

Linearity: The relationship between the independent variables and the dependent variable is linear.

Independence: The residuals (errors) are independent of each other.

Homoscedasticity: The variance of the residuals is constant across all levels of the independent variables.

Normality: The residuals are normally distributed.

No Multicollinearity: The independent variables should not be highly correlated with each other.

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)

yr	1069.78
september	370.72
light snow/rain	-454.45
spring	-391.91

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)

< Linear Regression is one of the simplest and most widely used statistical models. It is used to model the relationship between a dependent variable (target) and one or more independent variables (predictors). The goal of linear regression is to find the best-fit line (or hyperplane in the case of multiple features) that minimizes the difference between the predicted values and the actual values.

1. Overview of Linear Regression

In its simplest form (single variable or simple linear regression), the model predicts the dependent variable y as a linear function of one independent variable x :

$$y = \beta_0 + \beta_1 x + \epsilon$$

Where:

y is the dependent variable (the value you are trying to predict).

x is the independent variable (the feature you are using to predict y).

β_0 is the intercept (constant term).

β_1 is the coefficient or slope (the effect of the independent variable on the dependent variable).

ϵ is the error term or residual (the difference between the predicted and actual values).

In multiple linear regression, there are multiple independent variables, and the relationship becomes:

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p + \epsilon$$

Where:

x_1, x_2, \dots, x_p are the independent variables (features).

$\beta_1, \beta_2, \dots, \beta_p$ are the coefficients corresponding to each feature.

2. Objective of Linear Regression

The main objective of linear regression is to find the optimal values for the coefficients

$\beta_0, \beta_1, \dots, \beta_p$ such that the model minimizes the difference between the predicted values and the actual values. This is typically done using Least Squares Estimation (LSE).

3. Steps Involved in Linear Regression

Step 1: Model Representation

The linear regression model can be represented as:

$$Y = X\beta + \epsilon$$

Where:

Y is the vector of observed values of the dependent variable.

X is the matrix of independent variables (design matrix).

β is the vector of coefficients $[\beta_0, \beta_1, \dots, \beta_p]$.

ϵ is the error term (residuals).

Step 2: Least Squares Method

The primary goal is to minimize the sum of squared errors (residuals). The residual for each observation is the difference between the actual value and the predicted value:

$$\epsilon_i = y_i - \hat{y}_i$$

Where y_i is the actual value and \hat{y}_i is the predicted value.

To find the best-fitting line, we minimize the sum of squared residuals (RSS):

$$RSS = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

The solution to this minimization problem (i.e., finding the coefficients that minimize the residual sum of squares) can be derived using the Normal Equation:

$$\hat{\beta} = (X^T X)^{-1} X^T Y$$

Where:

$\hat{\beta}$ is the vector of estimated coefficients.

$X^T X$ is the transpose of the design matrix.

$(X^T X)^{-1}$ is the inverse of the matrix $X^T X$.

This equation gives the optimal values of the coefficients that minimize the sum of squared residuals.

Step 3: Model Fitting

After finding the optimal coefficients, the model is fit to the data, and predictions can be made.

For a new observation with features x_1, x_2, \dots, x_p , the predicted value is:

$$\hat{y} = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p$$

4. Assumptions of Linear Regression

For the model to produce reliable results, several assumptions must hold:

Linearity: The relationship between the dependent variable and the independent variables is linear.

This is the fundamental assumption of linear regression, which is why we use a straight line (or hyperplane in the case of multiple predictors).

Independence of Errors: The residuals (errors) should be independent of each other.

This assumption is critical for the statistical tests used to evaluate the model.

Homoscedasticity: The variance of the residuals should be constant across all levels of the independent variables.

If the variance of the errors is not constant (heteroscedasticity), it can lead to inefficient estimates of the coefficients.

Normality of Errors: The residuals should follow a normal distribution.

This assumption is important for hypothesis testing and confidence intervals of the regression coefficients.

No Multicollinearity: The independent variables should not be highly correlated with each other.

Multicollinearity can make the estimates of the regression coefficients unstable.

5. Model Evaluation

After fitting the model, several metrics are used to evaluate the performance of the linear regression model:

R-squared: This is the proportion of the variance in the dependent variable that is explained by the independent variables. It ranges from 0 to 1, with 1 indicating that the model explains all the variance.

$$R^2 = 1 - \frac{SS_{\text{residual}}}{SS_{\text{total}}} \quad R^2 = 1 - \frac{SS_{\text{total}}}{SS_{\text{residual}}}$$

Where SS_{residual} is the sum of squared residuals, and SS_{total} is the total sum of squares (variance of the target).

Mean Squared Error (MSE): This measures the average squared difference between the actual and predicted values. The smaller the MSE, the better the model.

$$MSE = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

Root Mean Squared Error (RMSE): This is the square root of MSE and represents the average error in the units of the dependent variable.

p-value: In hypothesis testing, p-values for each coefficient are used to determine if the corresponding feature significantly contributes to the model.

6. Advantages and Disadvantages of Linear Regression

Advantages:

Simple to implement: Linear regression is one of the easiest models to understand and implement.

Interpretability: The model coefficients have straightforward interpretations.

Efficiency: Linear regression is computationally efficient and fast to train on large datasets.

Disadvantages:

Assumptions: Linear regression makes several assumptions (linearity, normality, independence, etc.) that might not hold in real-world data, potentially leading to biased or incorrect results.

Sensitive to outliers: Linear regression is sensitive to outliers, which can disproportionately affect the model's predictions and coefficients.

Limited to linear relationships: If the true relationship between the features and target is non-linear, linear regression will not perform well.

7. Extensions of Linear Regression

Multiple Linear Regression: When there are more than one independent variable, linear regression extends to multiple linear regression.

Ridge and Lasso Regression: These are regularization techniques to prevent overfitting by adding a penalty term to the cost function.

Polynomial Regression: This extends linear regression by including polynomial terms (e.g., x^2, x^3) to model non-linear relationships.

>

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)

< Anscombe's Quartet is a set of four datasets created by the statistician **Francis Anscombe** in 1973. These datasets are famous for demonstrating how datasets with nearly identical statistical properties (e.g., mean, variance, correlation, and regression line) can exhibit vastly different

patterns when visualized.

The quartet underscores the importance of visualizing data rather than relying solely on summary statistics.

Properties of Anscombe's Quartet

Each of the four datasets in Anscombe's Quartet has the following nearly identical statistical properties:

Mean of x and y:

Mean of $x=9$ \text{Mean of } x = 9

Mean of $y=7.50$ \text{Mean of } y = 7.50

Variance of x and y:

Variance of $x=11$ \text{Variance of } x = 11

Variance of $y=4.12$ \text{Variance of } y = 4.12

Correlation between x and y:

$r=0.816$

Regression line:

The regression equation for all datasets is approximately: $y=3+0.5x$

Despite these similar numerical summaries, when the datasets are plotted, their differences become clear.

The Four Datasets in Anscombe's Quartet

Dataset 1:

A typical dataset where the points form a linear relationship.

The data aligns closely with the regression line.

Dataset 2:

A dataset with a perfect quadratic relationship.

While the linear regression equation is the same, the data clearly follows a curve.

Dataset 3:

A dataset where all data points except one form a perfect linear relationship.

A single outlier significantly affects the regression line and correlation.

Dataset 4:

A dataset where almost all points have the same x-value except for one.

The regression line is heavily influenced by this single influential data point.

Lessons from Anscombe's Quartet

The Importance of Visualization:

Summary statistics like mean, variance, and correlation can be misleading.

Visualization reveals the underlying patterns, relationships, and anomalies in the data.

Identifying Outliers and Influential Points:

Outliers and influential points can disproportionately affect regression lines and correlation coefficients.

Plotting the data can help identify these points.

Context in Statistical Analysis:

Data analysis must consider the context and structure of the data rather than relying solely on numbers.

Understanding the real-world implications of patterns is crucial.

Misleading Models:

The regression line might fit well statistically but fail to capture the true relationship in the data.

>

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)

< The Pearson's R (also known as Pearson's correlation coefficients) measures the strength between the different variables and the relation with each other. The Pearson's R returns values between -1 and 1. The interpretation of the coefficients are:

-1 coefficient indicates strong inversely proportional relationship.

0 coefficient indicates no relationship.

1 coefficient indicates strong proportional relationship.

$$r = \frac{n(\sum x * y) - (\sum x) * (\sum y)}{\sqrt{[n\sum x^2 - (\sum x)^2] * [n\sum y^2 - (\sum y)^2]}}$$

Where:

N = the number of pairs of scores

$\sum xy$ = the sum of the products of paired scores

$\sum x$ = the sum of x scores

$\sum y$ = the sum of y scores

$\sum x^2$ = the sum of squared x scores

$\sum y^2$ = the sum of squared y scores

>

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)

< What - The scaling is the data preparation step for regression model. The scaling normalizes these varied datatypes to a particular data range.

Why – Most of the times the feature data is collected at public domains where the interpretation of variables and units of those variables are kept open collect as much as possible. This results in to the high variance in units and ranges of data. If scaling is not done on these data sets, then the chances of processing the data without the appropriate unit conversion are high. Also the higher the range then higher the possibility that the coefficients are impaired to compare the dependent variable variance.

The scaling only affects the coefficients. The prediction and precision of prediction stays unaffected after scaling.

Normalization/Min-Max scaling – The Min max scaling normalizes the data within the range of 0 and 1. The Min max scaling helps to normalize the outliers as well.

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

Standardization converges all the data points into a standard normal distribution where mean is 0 and standard deviation is 1.

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

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)

<
$$VIF = \frac{1}{1-R^2}$$

The VIF formula clearly signifies when the VIF will be infinite. If the R² is 1 then the VIF is infinite. The reason for R² to be 1 is that there is a perfect correlation between 2 independent variables.

>

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)

< Q-Q plots are the quantile-quantile plots. It is a graphical tool to assess the 2 data sets are from common distribution. The theoretical distributions could be of type normal, exponential or uniform. The Q-Q plots are useful in the linear regression to identify the train data set and test data set are from the populations with same distributions. This is another method to check the normal distribution of the data sets in a straight line with patterns explained below

Interpretations

Similar distribution: If all the data points of quantile are lying around the straight line at an angle of 45 degree from x-axis.

Y values < X values: If y-values quantiles are lower than x-values quantiles.

X values < Y values: If x-values quantiles are lower than y-values quantiles.

Different distributions – If all the data points are lying away from the straight line.

Advantages

Distribution aspects like loc, scale shifts, symmetry changes and the outliers all can be identified from the single plot.

The plot has a provision to mention the sample size as well >
