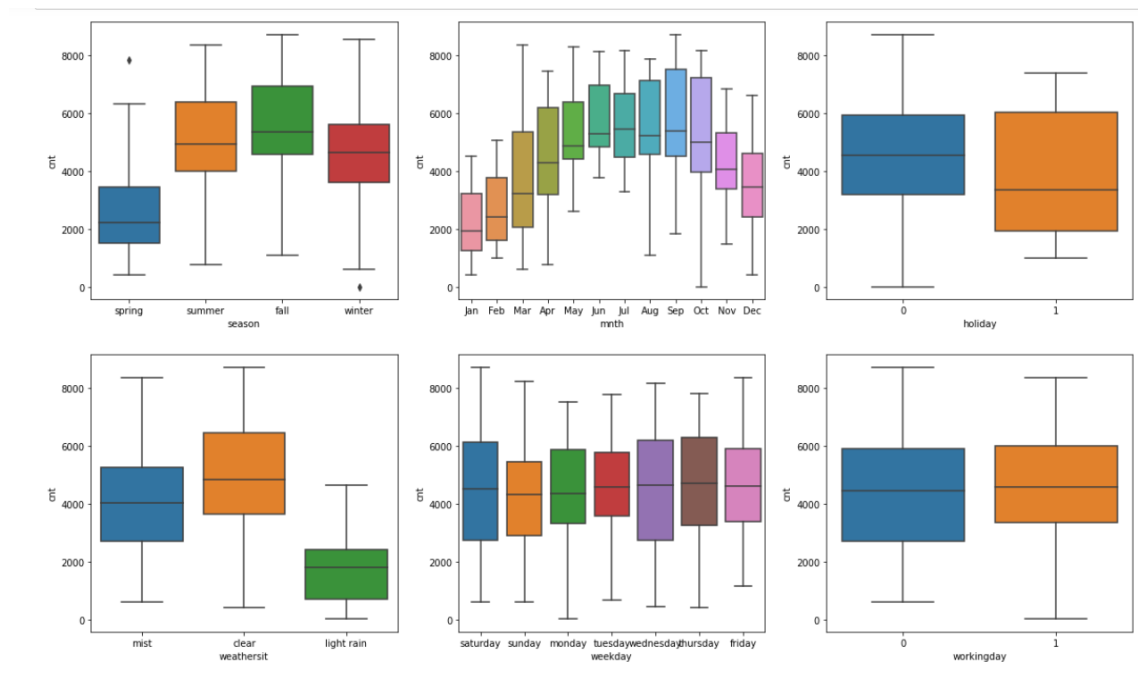


# Upgrad Assignment Questions

## 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?

Ans. From my analysis of the categorical variables from the dataset, I inferred that the during the fall and the summer season the registrations of the bikes were high. The people generally preferred to register and also book those bikes from the month of April till October as the numbers are high in those seasons. As you can see that the numbers indicate the preference of people to register those bikes in the clear weather during the summer and fall seasons.



2. Why is it important to use drop\_first=True during dummy variable creation?

Ans. When creating the dummy variables, it is important to use drop\_first=True because if we don't the column from which the dummy variables are being created then it can create the multicollinearity issues and it can be used to improve model interpretability and it also helps a great deal in improving the model efficiency and stability of the regression model.

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

Ans. Looking the pair-plot we can understand that **temp i.e. Temperature and atemp i.e. Feeling temperature** have the highest correlation with the target variable.

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

Ans. I validated the assumptions of the Linear regression by doing the following:

- Performing the Residual analysis to visualize and confirm that the residuals are normally distributed.
- Visualise the model to check for the liner relationship
- Calculate the VIF to ensure that there's a little or no multicollinearity.
- Performed Durbin-Watson test to check for 'No Autocorrelation' and the assumption is satisfied.
- Plotting the residuals for the homoscedasticity test.

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

Ans. Based on the final model, the top 3 features contributing significantly towards explaining the demand of the shared bikes are:

- Temperature
- Year
- Light Rain

## General Subjective Questions

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

Ans. Linear regression is a fundamental supervised machine learning algorithm used for modeling the relationship between a dependent variable (target) and one or more independent variables (features). It assumes a linear relationship between the input features and the target variable. The primary goal of linear regression is to find the best-fitting straight line that minimizes the difference between the predicted values and the actual target values.

The equation used in linear regression is :-  $y = mx + c$ , where  $y$  – dependent variable and  $x$  is the independent variable.

Here's a step-by-step breakdown of the linear regression algorithm:

1. Data Collection:

We need to collect a dataset on which the linear regression needs to be performed. The dataset needs to have the dependant variable and the independent variable. Each data point consists of a set of features and the associated target value.

## 2. Data Preprocessing:

Clean the data by handling missing values, outliers, and any other data quality issues. Normalize or standardize the features if necessary to ensure they're on similar scales.

## 3. Model Representation:

Linear regression assumes a linear relationship between the input features and the target. In the case of a single feature, the relationship can be represented as:

$$y = mx + b$$

Where:

- `y` is the target variable,
- `x` is the input feature,
- `m` is the slope (weight) of the line,
- `b` is the y-intercept (bias).

For multiple features, the equation becomes a linear combination of features with their corresponding weights and a bias term.

## 4. Cost Function:

The cost function (also known as the loss function) measures the difference between the predicted values and the actual target values. The goal is to minimize this difference. The most common cost function for linear regression is the Mean Squared Error (MSE).

$$MSE = (1/n) * \sum (y_{actual} - y_{predicted})^2$$

Where:

- `n` is the number of data points,
- `y\_actual` is the actual target value,
- `y\_predicted` is the predicted target value.

## 5. Gradient Descent:

Gradient descent is an optimization technique used to minimize the cost function. The algorithm iteratively adjusts the model's parameters (weights and biases) in the

direction that reduces the cost. The magnitude of the adjustments is determined by the learning rate, a hyperparameter.

#### 6. Parameter Updates:

During each iteration of gradient descent, the weights and biases are updated using the partial derivatives of the cost function with respect to the model's parameters. This update rule nudges the parameters towards values that minimize the cost.

#### 7. Training:

Repeatedly apply the gradient descent process on the training data until the algorithm converges to a point where the cost function is minimized or a stopping criterion is met.

#### 8. Prediction:

Once the model is trained and the optimal parameters are determined, you can use the model to make predictions on new, unseen data by plugging in the features and calculating the corresponding target value using the learned parameters.

In summary, linear regression aims to find the best-fitting linear relationship between input features and the target variable by minimizing the difference between predicted and actual values using gradient descent optimization. After this we need to perform the residual analysis and ensure that the residual errors are normally distributed with the mean as 0 and then the model evaluation process takes place.

#### 2. Explain the Anscombe's quartet in detail.

**Ans.** Anscombe's quartet is a set of four small datasets that have nearly identical statistical properties but exhibit very different visual patterns when plotted. These datasets were introduced by the statistician Francis Anscombe in 1973 to emphasize the importance of graphing data before performing any statistical analysis and to highlight the limitations of relying solely on summary statistics.

Each dataset in Anscombe's quartet consists of 11 (x, y) pairs of values, resulting in 11 data points. Here are the details of the four datasets:

#### 1. Dataset I:

- x: 10, 8, 13, 9, 11, 14, 6, 4, 12, 7, 5

- y: 8.04, 6.95, 7.58, 8.81, 8.33, 9.96, 7.24, 4.26, 10.84, 4.82, 5.68

#### 2. Dataset II:

- x: 10, 8, 13, 9, 11, 14, 6, 4, 12, 7, 5

- y: 9.14, 8.14, 8.74, 8.77, 9.26, 8.10, 6.13, 3.10, 9.13, 7.26, 4.74

#### 3. Dataset III:

- x: 10, 8, 13, 9, 11, 14, 6, 4, 12, 7, 5

- y: 7.46, 6.77, 12.74, 7.11, 7.81, 8.84, 6.08, 5.39, 8.15, 6.42, 5.73

#### 4. Dataset IV:

- x: 8, 8, 8, 8, 8, 8, 8, 19, 8, 8, 8

- y: 6.58, 5.76, 7.71, 8.84, 8.47, 7.04, 5.25, 12.50, 5.56, 7.91, 6.89

Despite having almost identical summary statistics (mean, variance, correlation, and linear regression parameters), these datasets showcase the danger of relying solely on numerical summaries to understand data. The visualizations of the datasets tell a different story:

1. Dataset I: A linear relationship, well-suited for linear regression.
2. Dataset II: A linear relationship with one outlier that influences the linear fit.
3. Dataset III: A nonlinear relationship that doesn't fit well with linear regression.
4. Dataset IV: A single outlier heavily influences the linear fit.

The main takeaway from Anscombe's quartet is that data visualization is essential for understanding the underlying patterns, relationships, and potential outliers in the data. It serves as a cautionary example against blindly applying statistical methods without first exploring the data graphically.

### 3. What is Pearson's R?

Ans. Pearson's correlation coefficient, often denoted as "r" or Pearson's "r," is a statistical measure that quantifies the strength and direction of the linear relationship between two continuous variables. It's a widely used method to assess how closely the points in a scatter plot lie along a straight line.

Pearson's correlation coefficient ranges from -1 to +1:

- A value of +1 indicates a perfect positive linear correlation, where as one variable increases, the other also increases proportionally.
- A value of -1 indicates a perfect negative linear correlation, where as one variable increases, the other decreases proportionally.
- A value close to 0 indicates a weak or no linear correlation between the variables.

Pearson's correlation coefficient measures the extent to which the points in a scatter plot cluster around the best-fit line. If the value of "r" is close to +1 or -1, it indicates a strong linear relationship. If "r" is close to 0, it suggests a weak or no linear correlation.

However, it's important to note that Pearson's correlation coefficient specifically measures linear relationships. It may not capture other types of relationships, such as nonlinear or monotonic relationships. Additionally, correlation does not imply causation; a high correlation does not necessarily mean one variable causes the other to change.

Pearson's correlation coefficient is commonly used in various fields, including statistics, data analysis, economics, social sciences, and more, to quantify and describe relationships between variables.

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

Ans. Scaling, in the context of data preprocessing for machine learning, refers to the process of transforming the features of a dataset to be on a similar scale. This is important because many machine learning algorithms are sensitive to the scale of input features. Scaling ensures that all features contribute equally to the model's training and that the algorithm performs better and converges faster.

Why Scaling is Performed:

Scaling is performed for several reasons:

1. **Algorithm Sensitivity:** Some machine learning algorithms, like k-nearest neighbors and gradient descent-based algorithms, are distance-based or rely on numerical optimization. These algorithms can be influenced by the magnitude of features, causing some features to dominate others.

2. **Faster Convergence:**Scaling can help gradient descent algorithms converge faster, as the optimization process becomes more stable when features are on similar scales.

3. **Regularization:** Some regularization techniques, like L1 and L2 regularization, assume that features are centered around zero. Scaling ensures this assumption holds.

4. **Distance Metrics:**Scaling is crucial for distance-based algorithms (e.g., k-means clustering), as the computation of distances between data points relies on feature magnitudes.

**Normalized Scaling vs. Standardized Scaling:**

Both normalized scaling and standardized scaling are methods to scale features, but they have different approaches and purposes:

1. **Normalized Scaling (Min-Max Scaling):**

Normalized scaling rescales features to a specified range, usually between 0 and 1. The formula for normalized scaling is

$$X_{\text{scaled}} = (X - X_{\text{min}}) / (X_{\text{max}} - X_{\text{min}})$$

Here, `X` is the original feature value, `X\_min` is the minimum value of the feature, and `X\_max` is the maximum value of the feature. Normalization preserves the relative relationships between the data points.

2. **Standardized Scaling (Z-Score Scaling):**

Standardized scaling transforms features so that they have a mean of 0 and a standard deviation of 1. The formula for standardized scaling is:

$$X_{\text{scaled}} = (X - \mu) / \sigma$$

Here, `X` is the original feature value, `μ` is the mean of the feature, and `σ` is the standard deviation of the feature. Standardization ensures that features have zero mean and unit variance.

**Key Difference:**

The main difference between normalized scaling and standardized scaling is that normalized scaling transforms features to a specific range (usually 0 to 1), while standardized scaling centers features around zero with a standard deviation of 1. Standardization is more appropriate when the distribution of the feature is close to normal, while normalization might be preferred when the distribution is skewed or when we want to maintain the original range of the data.

Ultimately, the choice between normalized scaling and standardized scaling depends on the characteristics of your data and the requirements of your specific machine learning algorithm.

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

Ans. VIF stands for Variance Inflation Factor, which is a measure used to assess multicollinearity in regression analysis. Multicollinearity occurs when two or more independent variables in a regression model are highly correlated with each other. High multicollinearity can lead to unstable coefficient estimates and difficulty in interpreting the effects of individual variables.

A common reason for observing infinite VIF values is a perfect linear relationship between one or more independent variables in the dataset. When two or more variables are perfectly correlated (linearly dependent), their relationship can be expressed as an exact linear equation. In such cases, the VIF for one of the correlated variables becomes mathematically infinite.

For example, let's say you have three variables A, B, and C, and you can express variable C as a linear combination of A and B ( $C = a * A + b * B$ ). In this case, the VIF for variable C would be infinite because it's perfectly predictable using the other variables.

It's important to note that while an infinite VIF value is a clear indicator of a problem with multicollinearity, other high VIF values (even if not infinite) can also indicate strong multicollinearity and may lead to unreliable coefficient estimates in regression analysis. In practice, when dealing with multicollinearity, it's important to identify the source of the issue, potentially through techniques like feature selection, dimensionality reduction, or domain knowledge to decide which variables to keep or combine.

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

Ans. A Q-Q plot, short for quantile-quantile plot, is a graphical tool used to assess whether a given dataset follows a particular theoretical distribution, such as the normal distribution. It's a way to visually compare the quantiles of the dataset to the quantiles of the expected theoretical distribution. Q-Q plots are particularly useful for understanding the distributional characteristics of the data and identifying deviations from a theoretical distribution.

How to Interpret a Q-Q Plot:

In a Q-Q plot:

- The x-axis represents the theoretical quantiles of the expected distribution.
- The y-axis represents the quantiles of the actual data.

If the data follows the expected distribution, the points in the Q-Q plot will lie approximately along a straight line. Deviations from a straight line suggest deviations from the expected distribution.

Use and Importance in Linear Regression:

Q-Q plots are especially relevant in linear regression for a couple of reasons:



1. Assumption Checking: One of the key assumptions of linear regression is that the errors (residuals) should be normally distributed. Q-Q plots can help you assess whether the residuals from a linear regression model follow a normal distribution. If the residuals deviate significantly from a straight line in the Q-Q plot, it might indicate non-normality in the residuals, suggesting that the assumption of normality might be violated.

2. Outlier Detection: Q-Q plots can also help identify outliers in the data. Outliers can be seen as data points that deviate substantially from the straight line in the Q-Q plot. Outliers can affect the fit of a linear regression model and might need special consideration or investigation.

3. Model Validity: If the residuals in a Q-Q plot deviate systematically from the expected straight line, it might indicate that the linear regression model is not capturing some underlying patterns in the data. This can lead to potential model improvement by considering more complex models or identifying additional features that might explain the variability in the data.

In summary, Q-Q plots are important tools in linear regression and statistical analysis in general. They help you assess the normality of data, detect outliers, and check the validity of assumptions underlying linear regression models. If a Q-Q plot suggests issues with normality or other assumptions, you might need to consider adjustments to your analysis or explore more advanced modeling techniques.