

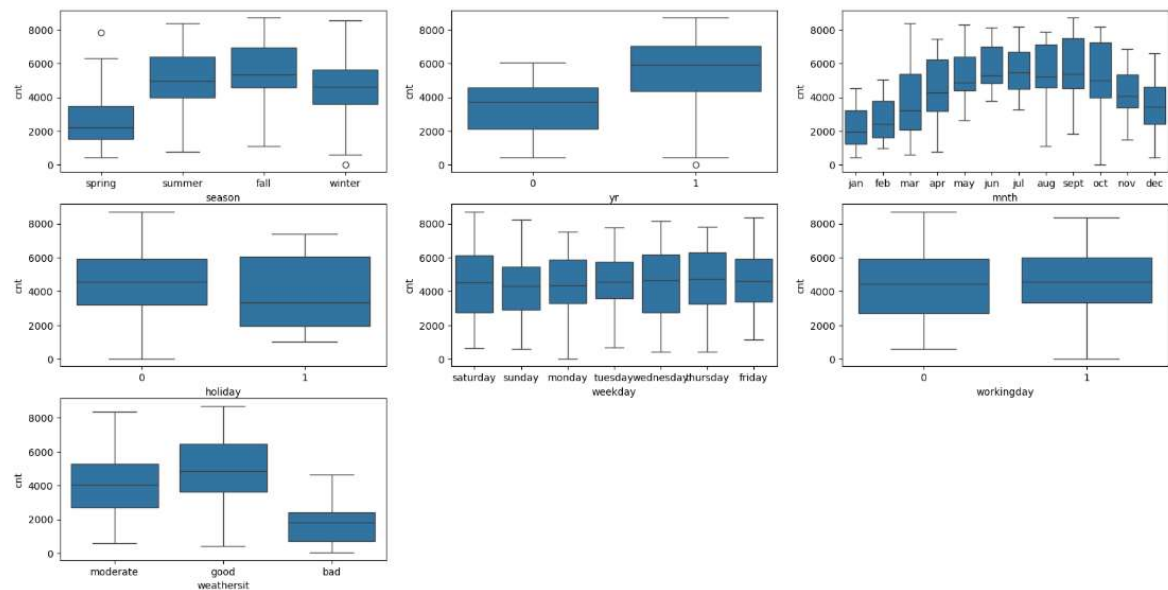
## Assignment-based Subjective Questions

**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)

The categorical variables like season, year, month, holiday, weekday, working day, weather sit, has major effect on the dependent variable 'cnt'. Below graphs shows the relation between 'cnt' & categorical variables,



- 
- Season - From season variable we can clearly say that Fall has the highest median, which means demand for bikes was highest during Fall & lowest during spring.
  - Yr - Year 2019 has higher demand compare to 2018
  - Month - Bikes were rented highest in September
  - Holiday - Rentals are less on holidays
  - Weekday - Bike demands are almost same for all days in a week
  - Workingday – There is no much difference in booking whether its working day or not
  - Weathersit – More bikes were rented when the weather is clear. Very less bikes were rented during bad weather condition.
- 

**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)

While creating dummy variables for a categorical variable you create 'n-1' new columns for 'n' values in the column, each indicating whether that level exists or not.

Using drop-first=True, the resultant can match up n-1 levels. Hence it reduces the correlation among the dummy variables.

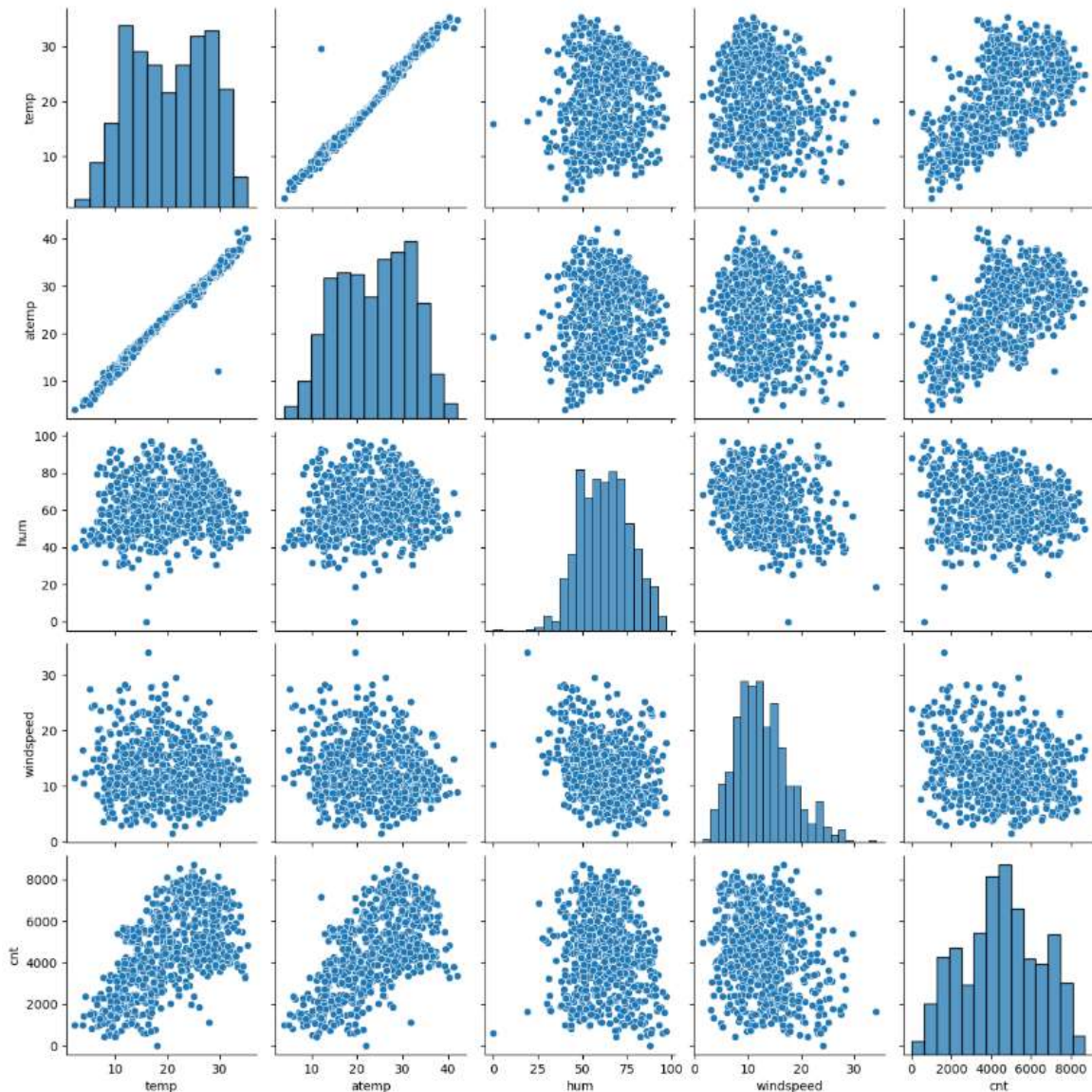
---

**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)

From the below figure we can clearly say that 'temp' & 'atemp' are highly correlated 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)

a. Linear Regression models are validated based on linearity between independent & dependent variables. Pairplots were used to visualize numeric variables & check if variables are linearly related or not.

b. Residual distribution should follow normal distribution & centered around mean = 0. We

validated this by performing residual analysis by using distplot of residuals.

c. We calculated VIF to quantify how strongly the feature variables are associated with one another. This handled multicollinearity in the data.

---

**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)

The top 3 significant features are, atemp (0.4290), yr(0.2419) & weathersit\_bad(-0.2742)

---

---

## 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 statistical method used to model the relationship between a dependent variable and one or more independent variables. It computes the linear relationship between the dependent variable and one or more independent features by fitting a linear equation with observed data. It predicts the continuous output variables based on the independent input variable.

The model's equation offers clear coefficients that illustrate the influence of each independent variable on the dependent variable, enhancing our understanding of the underlying relationships. Its simplicity is a significant advantage; linear regression is transparent, easy to implement, and serves as a foundational concept for more advanced algorithms.

The linear regression model gives a sloped straight line describing relationship within the variables.

A regression line can be positive or negative. The goal of the algorithm is to get the best values for  $a_0$  &  $a_1$  to find the best fitline & best fit line should have least error.

In linear regression, RFE or Mean squared error is used, which helps to figure out the best possible values for  $a_0$  &  $a_1$ .

---

**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)

<Your answer for Question 7 goes here>

Anscombe's quartet comprises a set of four datasets, having identical descriptive statistical properties in terms of means, variance, R-squared, correlations, and linear regression lines but having different representations when we scatter plots on a graph.

The datasets were created to demonstrate the importance of visualizing data and to show that summary statistics alone can be misleading.

The four datasets that make up Anscombe's quartet each include 11 x-y pairs of data. When plotted, each dataset seems to have a unique connection between x and y, with unique variability patterns and distinctive correlation strengths. Despite these variations, each dataset has the same summary statistics, such as the same x and y mean and variance, x and y correlation coefficient, and linear regression line

---

**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>

The Pearson correlation coefficient (r) is the most common way of measuring a linear correlation. It is a number between -1 and 1 that measures the strength and direction of the relationship between two variables.

Below is a formula for calculating the Pearson correlation coefficient (r):

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

Where

- N is sample size
  - X,y are the individual sample points indexed
- 

**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>

It is a step of data Pre-Processing which is applied to independent variables to normalize the data within a particular range. It also helps in speeding up the calculations in an algorithm. Most of the times, collected data set contains features highly varying in magnitudes, units and range. If scaling is not done then algorithm only takes magnitude in account and not units hence incorrect modelling. To solve this issue, we have to do scaling to bring all the variables to the same level of magnitude.

It is important to note that scaling just affects the coefficients and none of the other parameters like t-statistic, F-statistic, p-values, R-squared, etc.

**Normalization/Min-Max Scaling:**

---

- It brings all of the data in the range of 0 and 1.
  - `sklearn.preprocessing.MinMaxScaler` helps to implement normalization in python
- 

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

**Standardization Scaling:**

Standardization replaces the values by their Z scores. It brings all of the data into a standard normal distribution which has mean ( $\mu$ ) zero and standard deviation one ( $\sigma$ )

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

- 
- **sklearn.preprocessing.scale** helps to implement standardization in python.
  - One disadvantage of normalization over standardization is that it **loses** some information in the data, especially about **outliers**
- 

**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>

A VIF value becomes infinite when there is a perfect correlation between two or more independent variables in a regression model, meaning one variable can be perfectly predicted by a linear combination of the others, leading to a situation called "perfect multicollinearity" where the denominator in the VIF calculation becomes zero, resulting in an infinite value.

The value of VIF is calculated by the below formula:

$$VIF_i = \frac{1}{1 - R_i^2}$$

---

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

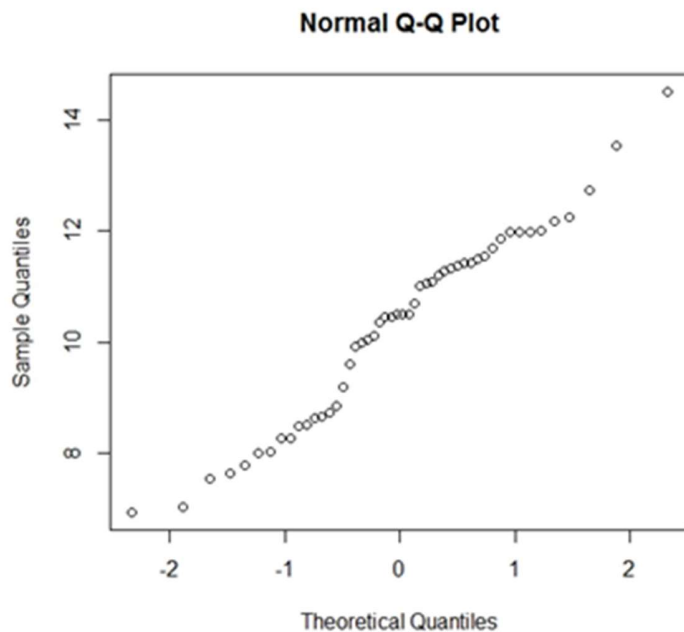
**Total Marks:** 3 marks (Do not edit)

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

<Your answer for Question 11 goes here>

The Q-Q plot or quantile-quantile plot is a graphical technique for determining if two data sets come from populations with a common distribution.

A Q-Q plot is a scatterplot created by plotting two sets of quantiles against one another. If both sets of quantiles came from the same distribution, we should see the points forming a line that's roughly straight. Here's an example of a Normal Q-Q plot when both sets of quantiles truly come from Normal distributions.



***Use of Q-Q plot in Linear Regression:***

The Q-Q plot is used to see if the points lie approximately on the line. If they don't, it means, our residuals aren't Gaussian (Normal) and thus, our errors are also not Gaussian.

***Importance of Q-Q plot: Below are the points:***

- I. The sample sizes do not need to be equal.
  - II. Many distributional aspects can be simultaneously tested. For example, shifts in location, shifts in scale, changes in symmetry, and the presence of outliers.
  - III. The q-q plot can provide more insight into the nature of the difference than analytical methods.
-