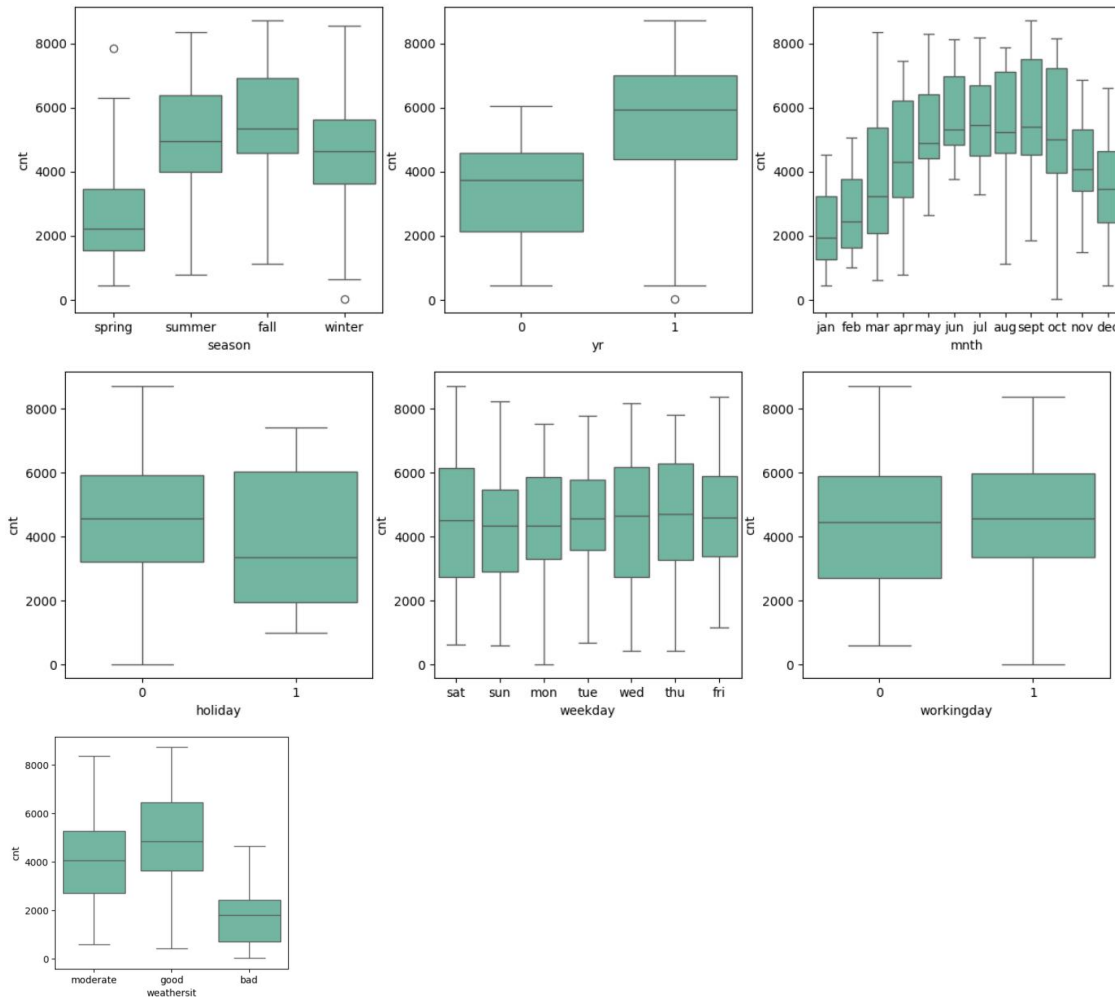


Assignment-based Subjective Questions

- From your analysis of the categorical variables from the dataset, what could you infer about their effect on the dependent variable?

Ans.



Categorical Variable	Inference
season	<ul style="list-style-type: none"> Fall, has the highest median, which shows that the demand was high during this season. It is least for spring.
yr(yr)	The year 2019 had a higher count of users as compared to the year 2018.
mnth(month)	<ul style="list-style-type: none"> There is gradual increase in no of rentals during Jan, Feb From Feb to March , the increase in rental count was steep and then increase gradually till June. June, July and August , almost have consistent demand for rentals September is the peak month for bike rentals October onwards rental count started declining till December.This observation is consistent with the advent of winter season in USA.
Holiday	Rentals reduced during holiday.
Weekday	The bike demand is almost constant throughout the week
Workingday	Median count of users is constant almost throughout the week. There is not much of difference in booking whether its working day or not.
Weathersit(Weather situation)	<ul style="list-style-type: none"> During bad weather (Heavy Rain + Ice Pallets + Thunderstorm + Mist, Snow + Fog) there are almost no users. Highest count was seen during good weather situation ie Clear, Few clouds,

Partly cloudy, Partly cloudy

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

Ans.

`drop_first=True` is important to use, as it helps in reducing the extra column created during dummy variable creation. Hence it reduces the correlations created among dummy variables.

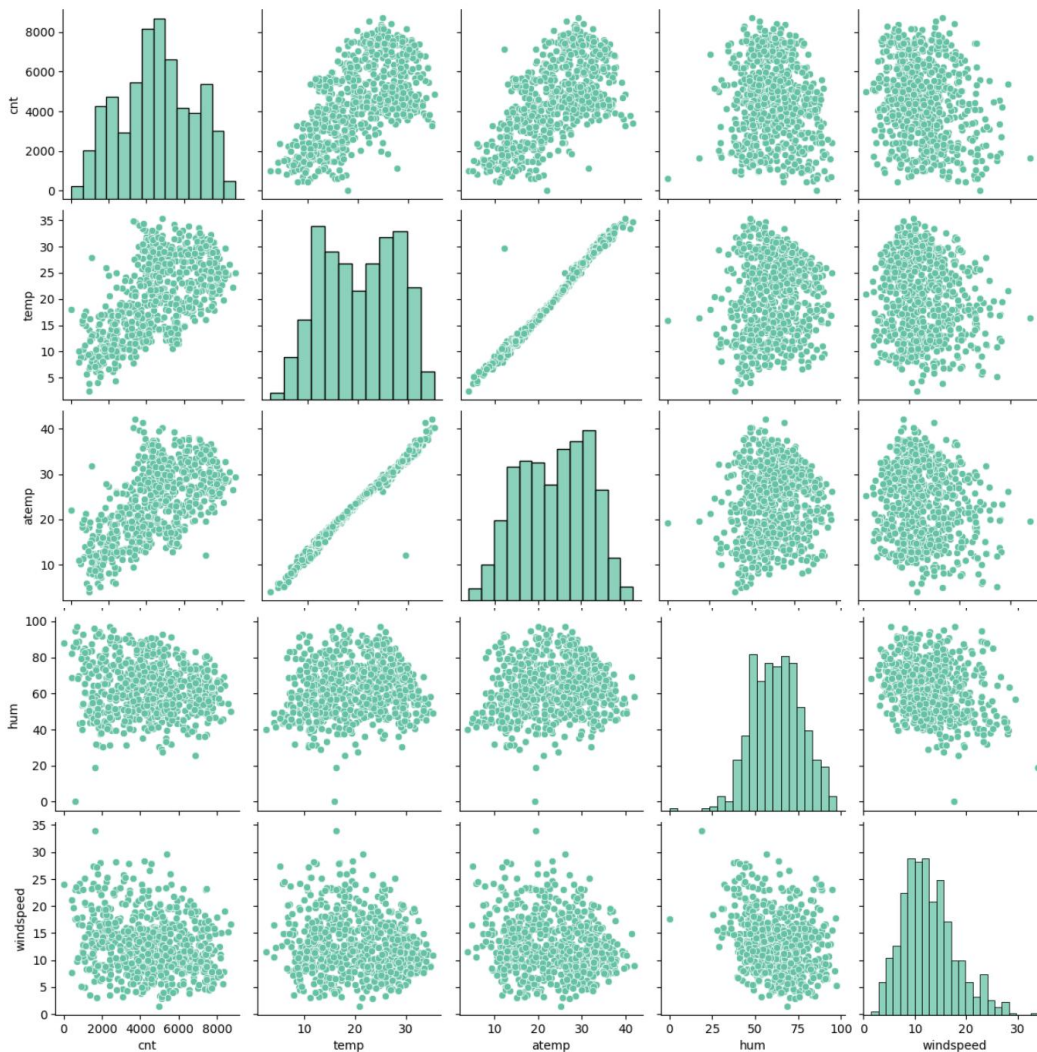
if we have categorical variable with n -levels, then we need to use $n-1$ columns to represent the dummy variables.

Consider a Categorical column with 3 types of values, we want to create dummy variable for that column. If one variable is neither furnished nor semi_furnished, then It is obvious unfurnished. So we do not need 3rd variable to identify the unfurnished.

Furnishing Status	Furnished	Unfurnished
Furnished	1	0
Semi- Furnished	0	1
Un-Furnished	0	0

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

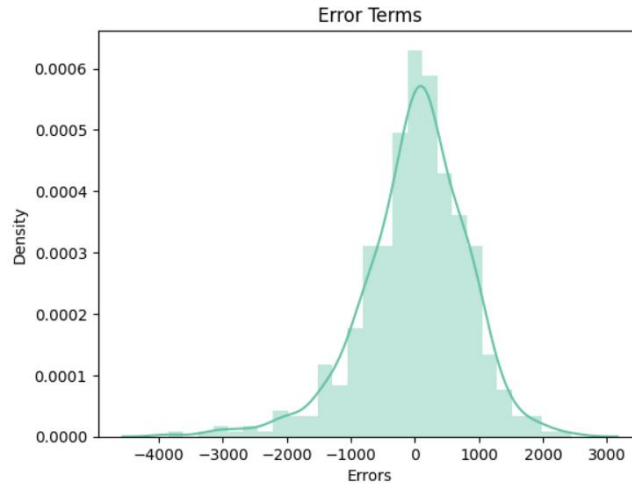
Ans. "temp" and "atemp" are the two numerical variables which are highly correlated with the target variable (cnt).



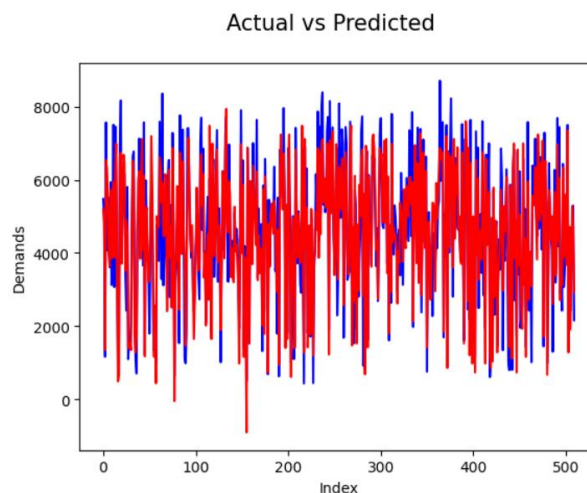
4. How did you validate the assumptions of Linear Regression after building the model on the training set?

Ans. We have done the following test to validate the assumption of Linear Regression:

- **Linearity Visualization** - We analyzed the pair plot to check linear relationship between independent and dependent variables. We plotted the numerical variables and found the linear relationship (Ref Question 3)
- **Residual Analysis** - Residuals distribution should follow normal distribution and centred around 0 (mean = 0). We validated this assumption by plotting a distplot of residuals and saw if residuals are following normal distribution or not.



- On training dataset, we plotted the actual and model predicted demand and check values almost followed the same pattern.



- Linear regression assumes that there is little or no multicollinearity in the data. Multicollinearity occurs when the independent variables are too highly correlated with each other. We calculated the VIF (Variance Inflation Factor) to quantify how strongly the feature variables in the new model are associated with one another

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

The top 3 significant features are:

1. temp - coefficient : 4048.1442
2. yr - coefficient : 1992.3065
3. weathersit_Light Snow & Rain & Scattered clouds - coefficient : -2313.5759

General Subjective Questions

1. Explain the linear regression algorithm in detail.

Ans.

Linear regression is a machine learning algorithm based on supervised learning. It performs a regression task, which means it predicts a continuous output variable (y) based on one or more input variables (x). It is mostly used for finding out the linear relationship between variables and forecasting.

The basic idea of linear regression is to find a line that best fits the data points, such that the distance between the line and the data points is minimized. The line can be represented by an equation of the form:

$$y = \theta_0 + \theta_1x$$

where θ_0 is the intercept (the value of y when x is zero) and θ_1 is the slope (the change in y for a unit change in x). These are called the parameters or coefficients of the linear model.

To find the best values of θ_0 and θ_1 , we need to define a cost function that measures how well the line fits the data. A common choice is the mean squared error (MSE), which is the average of the squared differences between the actual y values and the predicted y values:

$$MSE = (1/n) * \sum (y - y')^2$$

where n is the number of data points, y is the actual value, and y' is the predicted value.

The goal is to minimize the MSE by adjusting θ_0 and θ_1 . There are different methods to do this, such as gradient descent, normal equation, or using libraries like scikit-learn.

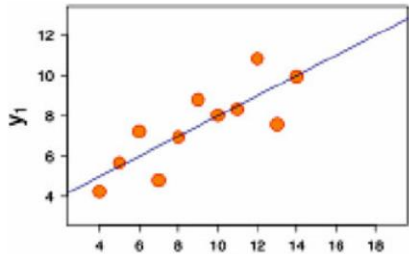
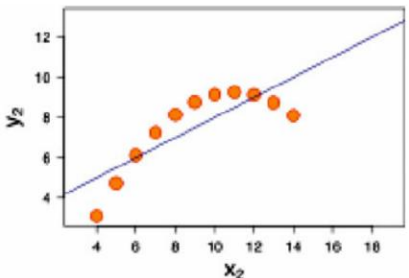
Linear regression can also be extended to multiple input variables (x_1, x_2, \dots, x_n), in which case the equation becomes:

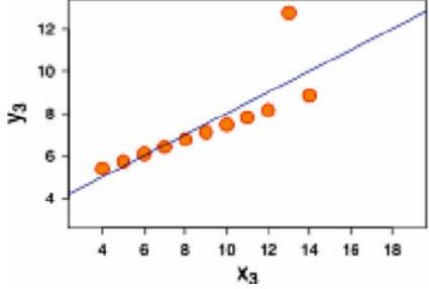
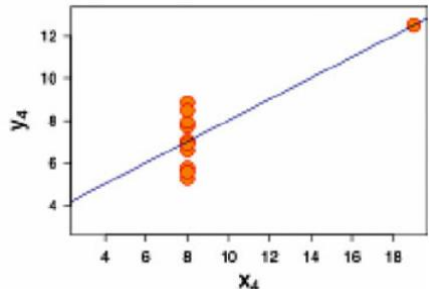
$$y = \theta_0 + \theta_1x_1 + \theta_2x_2 + \dots + \theta_nx_n$$

Limitations are: it assumes a linear relationship between the input variables and the output variable, which may not always be the case. Another limitation is that it may be sensitive to outliers or multicollinearity.

2. Explain the Anscombe’s quartet in detail.
 Ans.

Anscombe’s Quartet was developed by statistician Francis Anscombe. It includes four data sets that have almost identical statistical features, but they have a very different distribution and look totally different when plotted on a graph. It was developed to emphasize both the importance of graphing data before analyzing it and the effect of outliers and other influential observations on statistical properties.

	<p>The graph shows simple linear relationship.</p>
	<p>The graph is not distributed normally; the relation between them is not linear.</p>

	<p>In this graph, the distribution is linear, but should have a different regression line the calculated regression is offset by the one outlier which exerts enough influence to lower the correlation coefficient from 1 to 0.816.</p>
	<p>This graph shows an example when one high-leverage point is enough to produce a high correlation coefficient, even though the other data points do not indicate any relationship between the variables.</p>

3. What is Pearson's R?

Ans.

Pearson's r is a numerical summary of the strength of the linear association between the variables. Its value ranges between -1 to +1. It shows the linear relationship between two sets of data. In simple terms, it tells us "can we draw a line graph to represent the data?"

Formula

$$r = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2 \sum (y_i - \bar{y})^2}}$$

r = correlation coefficient

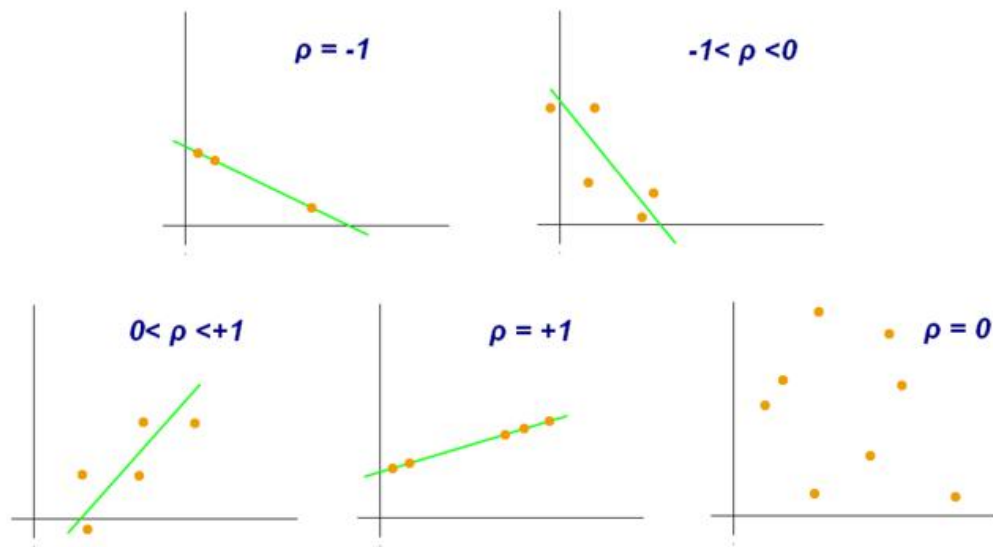
x_i = values of the x-variable in a sample

\bar{x} = mean of the values of the x-variable

y_i = values of the y-variable in a sample

\bar{y} = mean of the values of the y-variable

As can be seen from the graph below, $r = 1$ means the data is perfectly linear with a positive slope $r = -1$ means the data is perfectly linear with a negative slope $r = 0$ means there is no linear association



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

Ans.

Feature scaling is a method used to normalize or standardize the range of independent variables or features of data. It is performed during the data pre-processing stage to deal with varying values in the data set. If feature scaling is not done, then a machine learning algorithm tends to weigh greater values, higher and consider smaller values as the lower values, irrespective of the units of the values.

- Normalization is generally used when you know that the distribution of your data does not follow a Gaussian distribution. This can be useful in algorithms that do not assume any distribution of the data like K-Nearest Neighbours and Neural Networks.
- Standardization, on the other hand, can be helpful in cases where the data follows a Gaussian distribution. However, this does not have to be necessarily true. Also, unlike normalization, standardization does not have a bounding range. So, even if you have outliers in your data, they will not be affected by standardization.

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

Ans.

The VIF (Variance Inflation Factor) gives how much the variance of the coefficient estimate is being inflated by collinearity. If there is perfect correlation, then $VIF = \text{infinity}$. It gives a basic quantitative idea about how much the feature variables are correlated with each other. It is an extremely important parameter to test our linear model.

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

Where R^2 is the R-square value of that independent variable which we want to check how well this independent variable is explained well by other independent variables. If that independent variable can be explained perfectly by other independent variables, then it will have perfect correlation and its R^2 value will be equal to 1. So, $VIF = 1/(1-R^2)$ which gives $VIF = 1/0$ which results in "infinity". The numerical value for VIF tells you (in decimal form) what percentage the variance (i.e. the standard error squared) is inflated for each coefficient. For example, a VIF of 1.9 tells you that the variance of a particular coefficient is 90% bigger than what you would expect if there was no multicollinearity — if there was no correlation with other predictors.

A rule of thumb for interpreting the variance inflation factor:

- 1 = not correlated.
- Between 1 and 5 = moderately correlated.
- Greater than 5 = highly correlated.

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 is a plot of the quantiles of the first data set against the quantiles of the second data set. It is used to compare the shapes of distributions. 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.

The q-q plot is used to answer the following questions:

- Do two data sets come from populations with a common distribution?
- Do two data sets have common location and scale?
- Do two data sets have similar distributional shapes?
- Do two data sets have similar tail behaviour?

