1. From your analysis of the categorical variables from the dataset, what could you infer about their    effect on the dependent variable?
Ans: From year variable it shows that more bikes were rented in 2019.
     Season box plots indicates that more bikes are rent during fall season.
     Working day and holiday box plots indicate that more bikes are rent during normal working days        than on weekends or holidays.
     The month box plots indicates that more bikes are rent during Aug, sep, Oct month
     Weathersit box plots indicates that more bikes are rent during Clear, Few clouds, Partly cloudy      weather.

2. Why is it important to use drop_first=True during dummy variable creation?
Ans: It helps in reducing the extra column created during dummy variable creation. Hence it reduces the       correlations created among dummy variables.


3. Looking at the pair-plot among the numerical variables, which one has the highest correlation with the target variable?
Ans: Month, Season, temp, atemp, year, season, windspeed - These predictor variables have highest corr.


4. How did you validate the assumptions of Linear Regression after building the model on the training set?
Ans: I did the residual analysis to validate LR assumtions. Same can be seen in plots in notebook.


5. Based on the final model, which are the top 3 features contributing significantly towards explaining the demand of the shared bikes?
Ans: Season, year and temp/atemp


General Subjective Questions

1. Explain the linear regression algorithm in detail.
Linear regression is one of the very basic forms of machine learning where we train a model to predict the behaviour of your data based on some variables. In the case of linear regression as you can see the name suggests linear that means the two variables which are on the x-axis and y-axis should be linearly correlated.

Linear regression provides a powerful statistical method to find the relationship between variables. It hardly needs further tuning. However, it's only limited to linear relationships.

Linear regression produces the best predictive accuracy for linear relationship whereas its little sensitive to outliers and only looks at the mean of the dependent variable.


2. Explain the Anscombe's quartet in detail.
Anscombe's Quartet can be defined as a group of four data sets which are nearly identical in simple descriptive statistics, but there are some peculiarities in the dataset that fools the regression model if built. They have very different distributions and appear differently when plotted on scatter plots.

It was constructed in 1973 by statistician Francis Anscombe to illustrate the importance of plotting the graphs before analyzing and model building, and the effect of other observations on statistical properties.There are these four data set plots which have nearly same statistical observations, which provides same statistical information that involves variance, and mean of all x,y points in all four datasets.

3. What is Pearson's R?
Pearson's correlation coefficient is the covariance of the two variables divided by the product of their standard deviations. The form of the definition involves a "product moment", that is, the mean (the first moment about the origin) of the product of the mean-adjusted random variables; hence the modifier product-moment in the name.

Pearson's Correlation Coefficient is named after Karl Pearson. He formulated the correlation coefficient from a related idea by Francis Galton in the 1880s.
Using the formula proposed by Karl Pearson, we can calculate a linear relationship between the two given variables. For example, a child's height increases with his increasing age (different factors affect this biological change). So, we can calculate the relationship between these two variables by obtaining the value of Pearson's Correlation Coefficient r. There are certain requirements for Pearson's Correlation Coefficient:

> Scale of measurement should be interval or ratio

> Variables should be approximately normally distributed

> The association should be linear

> There should be no outliers in the data

4. What is scaling? Why is scaling performed? What is the difference between normalized scaling and standardized scaling?
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
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 ($\sigma$).

5. You might have observed that sometimes the value of VIF is infinite. Why does this happen?
An infinite VIF value indicates that the corresponding variable may be expressed exactly by a linear combination of other variables (which show an infinite VIF as well).

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

Q-Q Plots (Quantile-Quantile plots) are plots of two quantiles against each other. A quantile is a fraction where certain values fall below that quantile. For example, the median is a quantile where 50% of the data fall below that point and 50% lie above it. The purpose of Q Q plots is to find out if two sets of data come from the same distribution. A 45 degree angle is plotted on the Q Q plot; if the two data sets come from a common distribution, the points will fall on that reference line.