# 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?

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

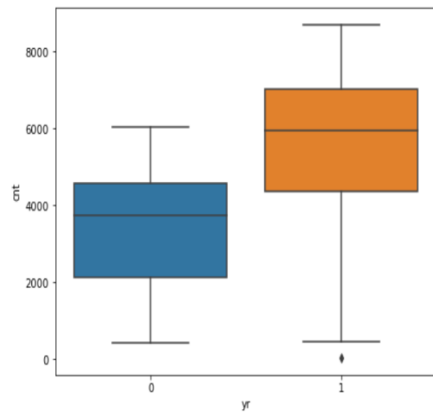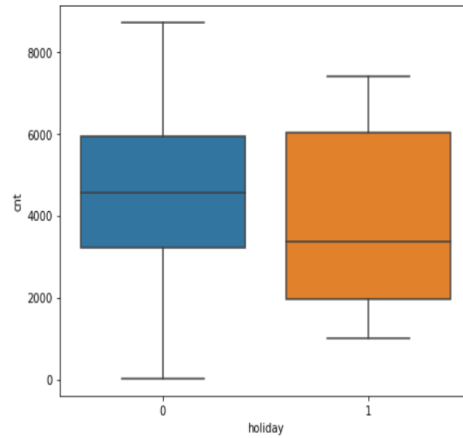   Answer: - drop_first=True is important to use, as it helps in reducing the extra column created during dummy variable creation. Hence it reduces the correlations created among dummy variables

   Suppose, categorical variable are n it's must have n-1 dummy variable

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

   Answer: - atemp.

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

Answer: -
1. There should be a linear and additive relationship between dependent variable and independent variables
2. There should be no correlation between the residual (error) terms.
3. The independent variables should not be correlated. Absence of this phenomenon is known as multicollinearity.
4. The error terms must have constant variance.
5.  The error terms must be normally distributed.

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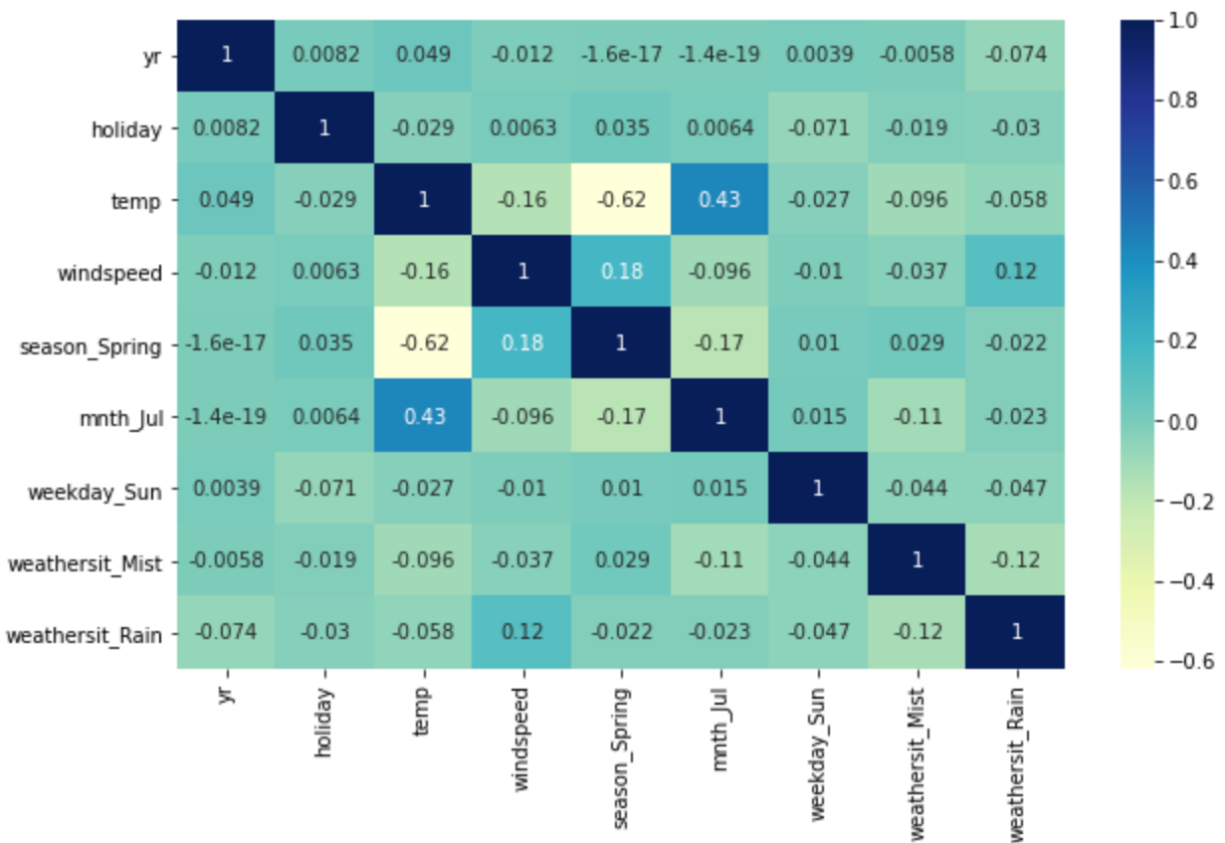
Answer: -

1. Temp

2. holiday,

3. weekday_sun

| | yr | holiday | temp | windspeed | season_Spring | mnth_Jul | weekday_Sun | weathersit_Mist | weathersit_Rain |
|---|---|---|---|---|---|---|---|---|---|
| yr | 1 | 0.0082 | 0.049 | -0.012 | -1.6e-17 | -1.4e-19 | 0.0039 | -0.0058 | -0.074 |
| holiday | 0.0082 | 1 | -0.029 | 0.0063 | 0.035 | 0.0064 | -0.071 | -0.019 | -0.03 |
| temp | 0.049 | -0.029 | 1 | -0.16 | -0.62 | 0.43 | -0.027 | -0.096 | -0.058 |
| windspeed | -0.012 | 0.0063 | -0.16 | 1 | 0.18 | -0.096 | -0.01 | -0.037 | 0.12 |
| season_Spring | -1.6e-17 | 0.035 | -0.62 | 0.18 | 1 | -0.17 | 0.01 | 0.029 | -0.022 |
| mnth_Jul | -1.4e-19 | 0.0064 | 0.43 | -0.096 | -0.17 | 1 | 0.015 | -0.11 | -0.023 |
| weekday_Sun | 0.0039 | -0.071 | -0.027 | -0.01 | 0.01 | 0.015 | 1 | -0.044 | -0.047 |
| weathersit_Mist | -0.0058 | -0.019 | -0.096 | -0.037 | 0.029 | -0.11 | -0.044 | 1 | -0.12 |
| weathersit_Rain | -0.074 | -0.03 | -0.058 | 0.12 | -0.022 | -0.023 | -0.047 | -0.12 | 1 |

# General Subjective Questions

1. Explain the linear regression algorithm in detail.

   Answer: - Linear Regression is a machine learning algorithm based on supervised learning. It performs a regression task. Regression models a target prediction value based on independent variables. It is mostly used for finding out the relationship between variables and forecasting. Different regression models differ based on – the kind of relationship between dependent and independent variables

2. Explain the Anscombe's quartet in detail.

   Answer: - Anscombe's quartet comprises four datasets that have nearly identical simple statistical properties yet appear very different when graphed. Each dataset consists of eleven (x,y) points. They were constructed in 1973 by the statistician Francis Anscombe to demonstrate both the importance of graphing data before analyzing it and the effect of outliers on statistical properties.

   The quartet is still often used to illustrate the importance of looking at a set of data graphically before starting to analyze according to a particular type of relationship, and the inadequacy of basic statistic properties for describing realistic datasets.

3. What is Pearson's R?

   Answer: - Pearson's r is a numerical summary of the strength of the linear association between the variables. If the variables tend to go up and down together, the correlation coefficient will be positive. If the variables tend to go up and down in opposition with low values of one variable associated with high values of the other, the correlation coefficient will be negative.

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

   Answer: - 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 must 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 the data in the range of 0 and 1.
   **Standardization Scaling:** Standardization replaces the values by their Z scores. It brings all the data into a standard normal distribution which has mean (**μ**) zero and standard deviation one (**σ**).

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

   Answer: - If there is perfect correlation, then VIF = infinity. This shows a perfect correlation between two independent variables. In the case of perfect correlation, we get R2 =1, which lead to 1/(1-R2) infinity. To solve this problem, we need to drop one of the variables from the dataset which is causing this perfect multicollinearity.
   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.

   Answer: - 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.

   A Q–Q plot is used to compare the shapes of distributions, providing a graphical view of how properties such as location, scale, and skewness are similar or different in the two distributions.