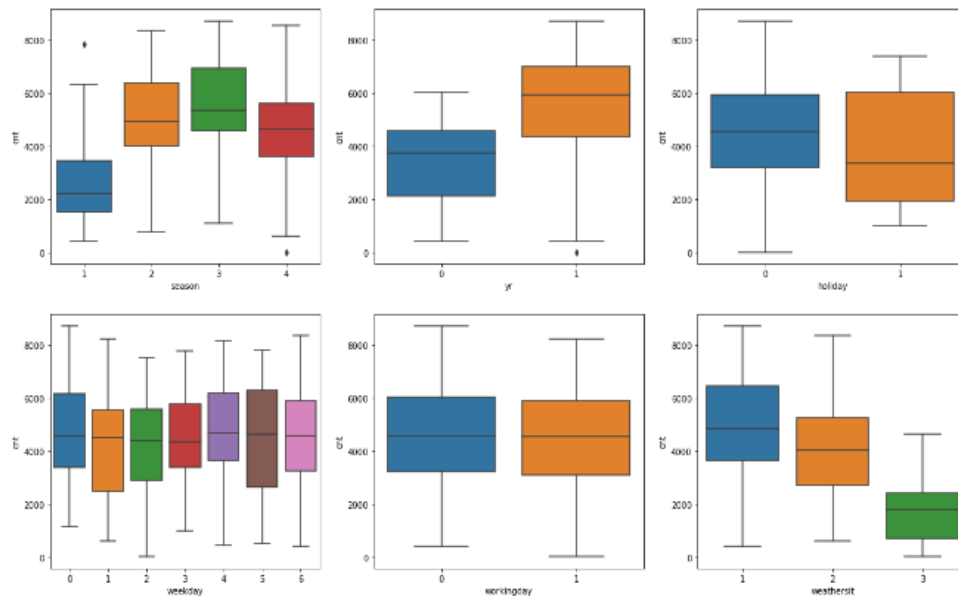# Assignment-based Subjective Questions

1. From your analysis of the categorical variables from the dataset, what could you infer about their effect on the dependent variable? (3 marks)

Ans. -



a) Demand for share bike was max in season 'FALL' and min in season 'SPRING'.
b) Demand for share bike in year 2019 was higher than demand in year 2018.
c) Demand for share bike was higher when weather was clear and there was no any demand when there was heavy rain.

2. Why is it important to use drop_first=True during dummy variable creation? (2 mark)
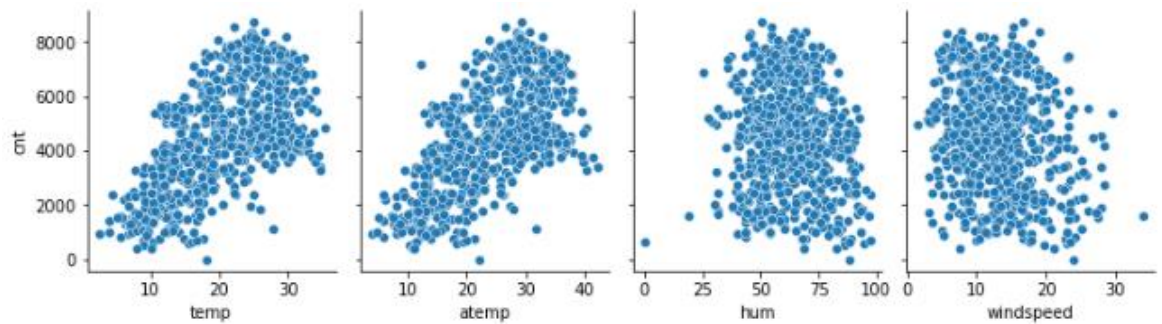
Ans. –

a) If there are 6 dummy variables then we can predict output of any variable by using remaining 5 variables. So we can drop any dummy variable. By doing this we dropped unnecessary column.

3. Looking at the pair-plot among the numerical variables, which one has the highest correlation with the target variable? (1 mark)
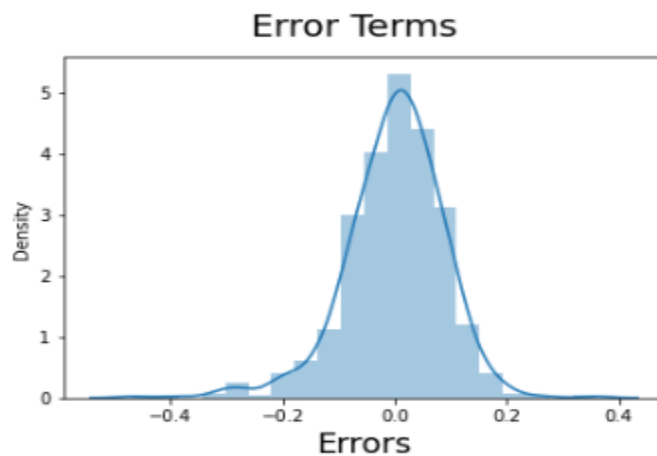
Ans. –

a) From below diagram we can say that Temperature has highest correlation with the Target variable .
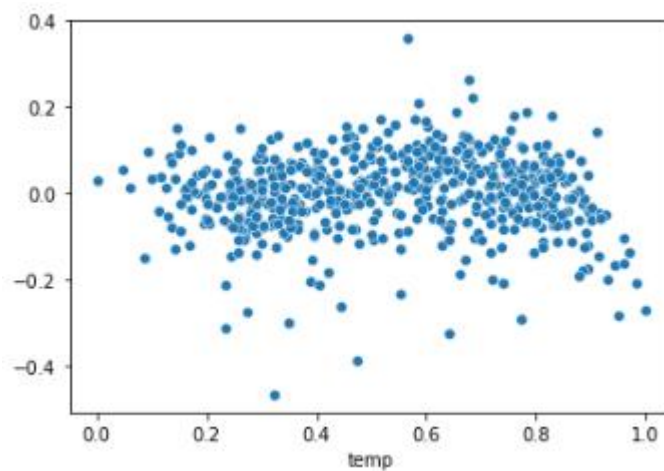b) And that is positive correlation.

4. How did you validate the assumptions of Linear Regression after building the model on the training set? (3 marks)

Ans.-

a) The assumption of normality is made, as it has been observed that the error terms generally follow a normal distribution with mean equal to zero in most cases.



b) The error terms should not be dependent on one another
Ex. – scatterplot of temp vs error terms



c) The variance should not increase (or decrease) as the error values change. Also, the variance should not follow any pattern as the error terms change.

5) Based on the final model, which are the top 3 features contributing significantly towards explaining the demand of the shared bikes? (2 marks)

Ans.-

a) Temperature (0.4108)
b) Light-snow (-0.2977)
c) Year (0.2331)

# General Subjective Questions

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

Ans.-

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 they are considering, and the number of independent variables getting used.

2. Explain the Anscombe's quartet in detail. (3 marks)

Ans.-

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? (3 marks)

Ans.-

Pearson's R is a measure of linear correlation between two sets of data. It is the ratio between the covariance of two variables and the product of their standard deviations; thus, it is essentially a normalized measurement of the covariance, such that the result always has a value between −1 and 1. As with covariance itself, the measure can only reflect a linear correlation of variables, and ignores many other types of relationships or correlations. As a simple example, one would expect the age and height of a sample of teenagers from a high school to have a Pearson correlation coefficient significantly greater than 0, but less than 1 (as 1 would represent an unrealistically perfect correlation).

4. What is scaling? Why is scaling performed? What is the difference between normalized scaling and standardized scaling? (3 marks)

Ans.-

**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.

**Why Scaling:** 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.

**Normalized :** It brings all of the data in the range of 0 and 1.

**Standardized:** 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? (3 marks)

Ans.-

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. (3 marks)

Ans.-

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. If the two distributions being compared are similar, the points in the Q–Q plot will approximately lie on the line $y = x$. If the distributions are linearly related, the points in the Q–Q plot will approximately lie on a line, but not necessarily on the line $y = x$. Q–Q plots can also be used as a graphical means of estimating parameters in a location-scale family of distributions.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.