

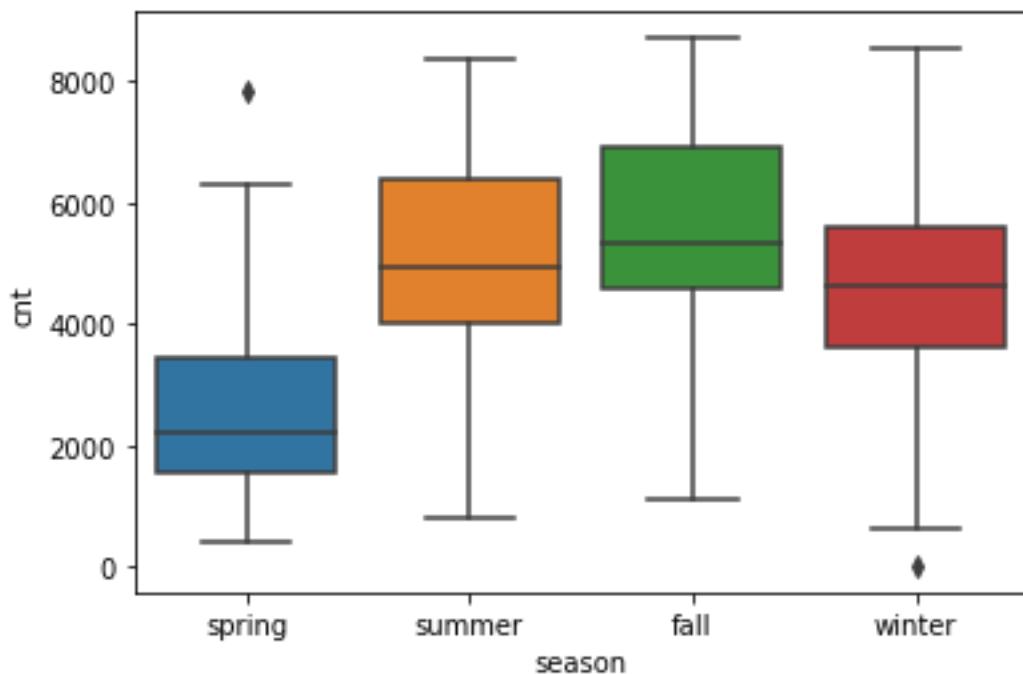
Assignment-based Subjective Questions

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

Ans:

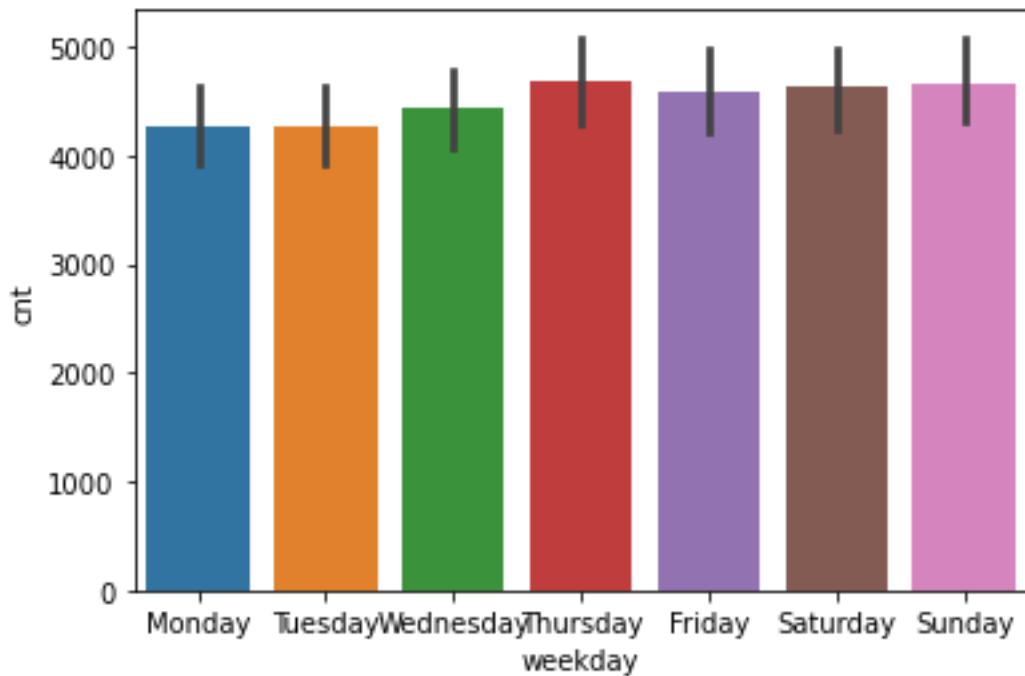
- Impact of categorical variable ‘season’:

Count of rental bikes is the highest in fall season followed by the Summer and the Winter. In the Spring time demand is the least.



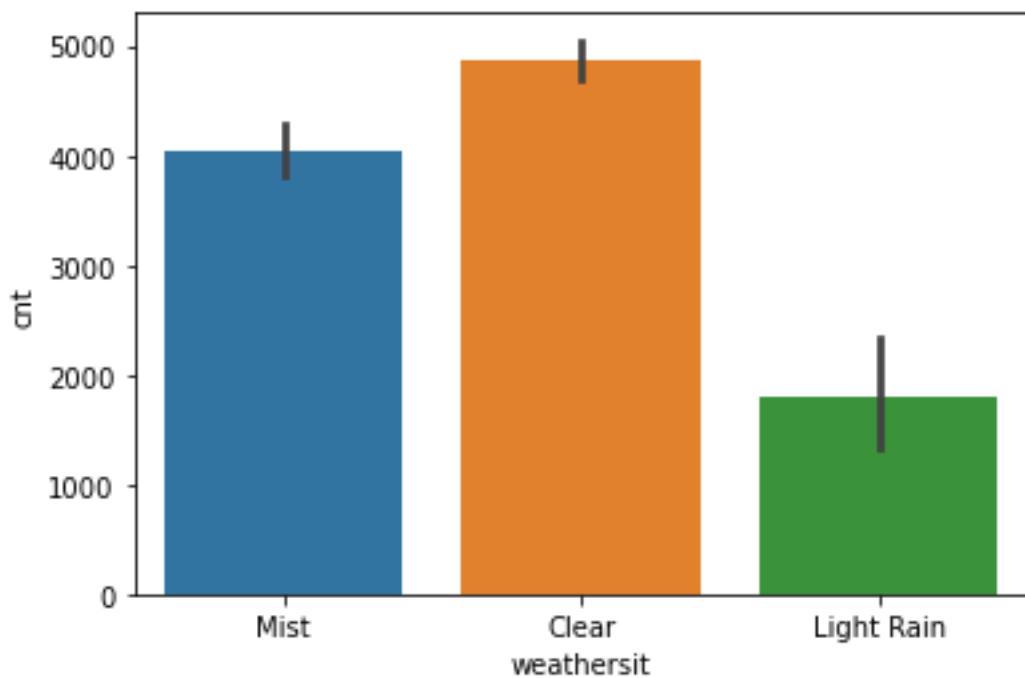
- Impact of categorical variable ‘weekday’:

The count of rental bikes is slightly higher on Thursday than other days. On Friday , Saturday and Sunday the count is almost same. Monday and Tuesday has little less demand. Avg. Demand of rental bike lies between 4000 to 4500.

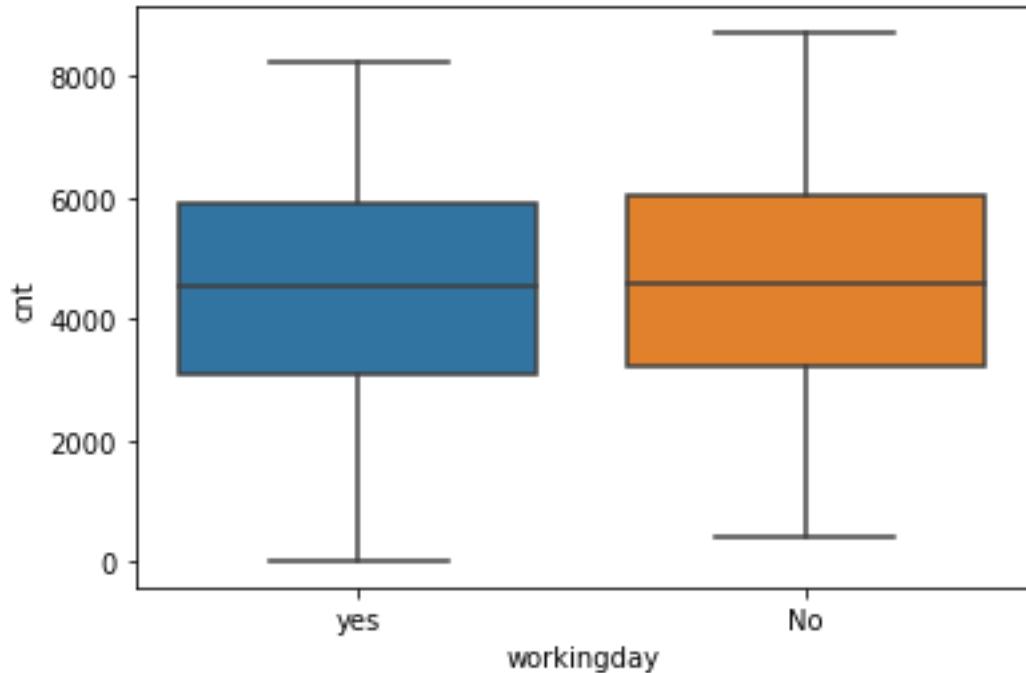


- Impact of categorical variable ‘weathersit’:

Count of rental bike goes highest in clear weather and decreases in rainy day.

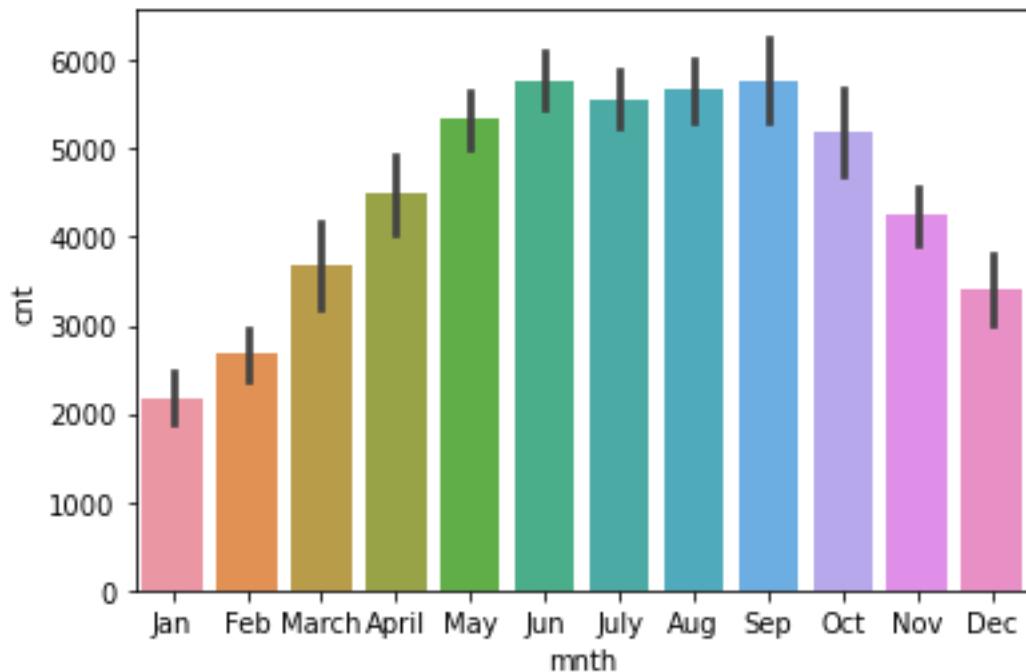


- Impact of categorical variable ‘workingday’:



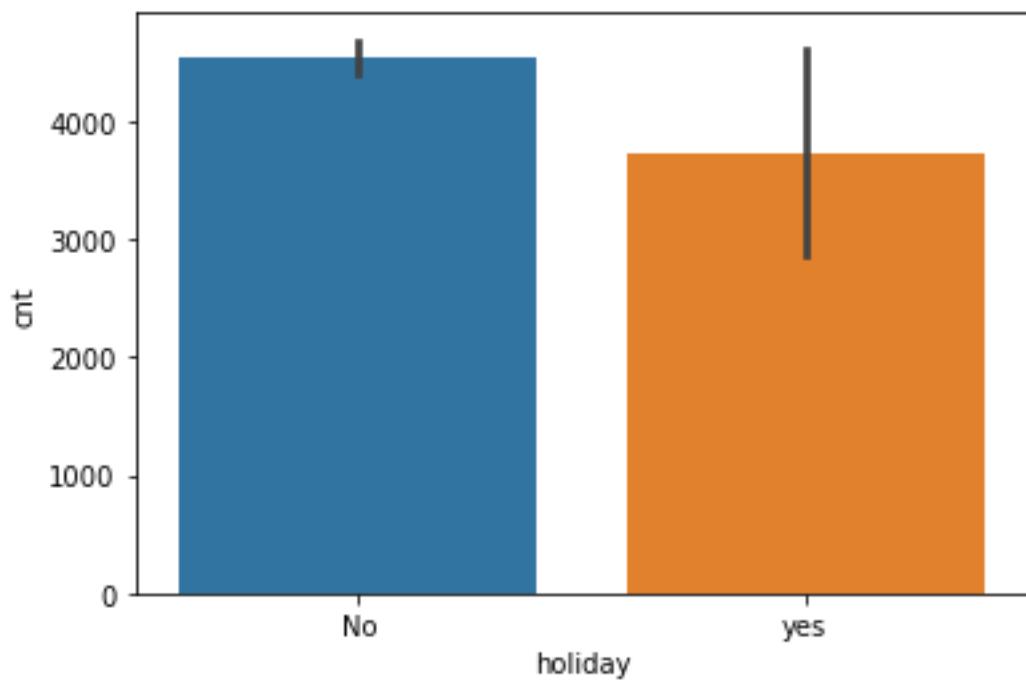
There is slightly high demand of shared bikes on non-working day

- Impact of categorical variable ‘mnth’:



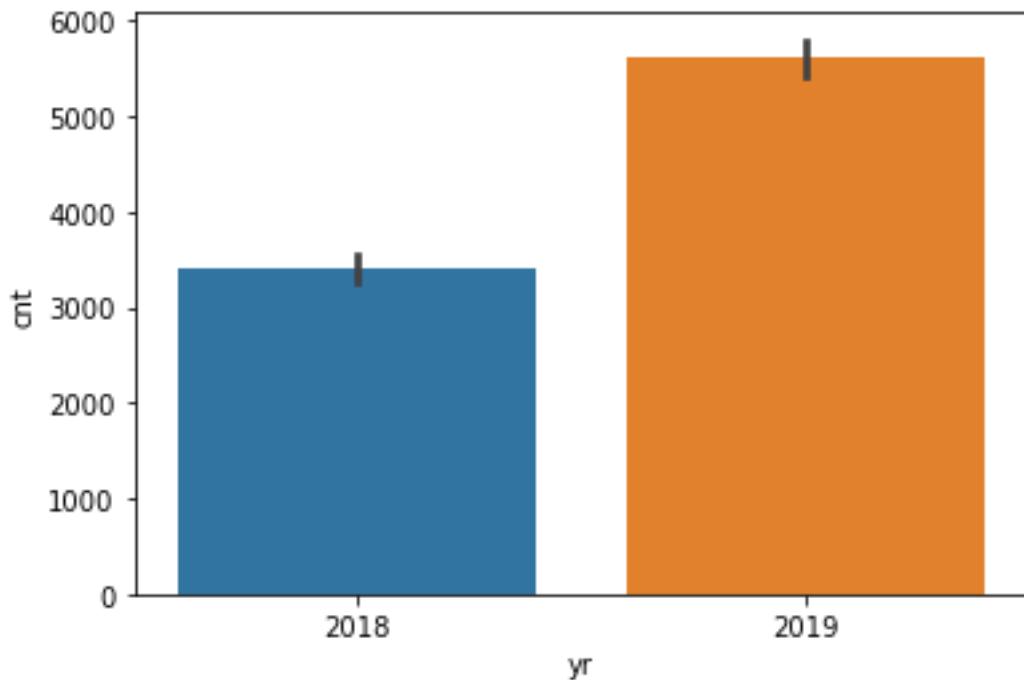
Demand straightly increases up to month of Jun then deteriorate but again rise is shown from July to August and then again deteriorate continuously till December with lowest value.

- Impact of categorical variable ‘holiday’:



Demand is high on non-holidays,

- Impact of categorical variable ‘yr’:



Demand increased sharply in the year 2019.

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

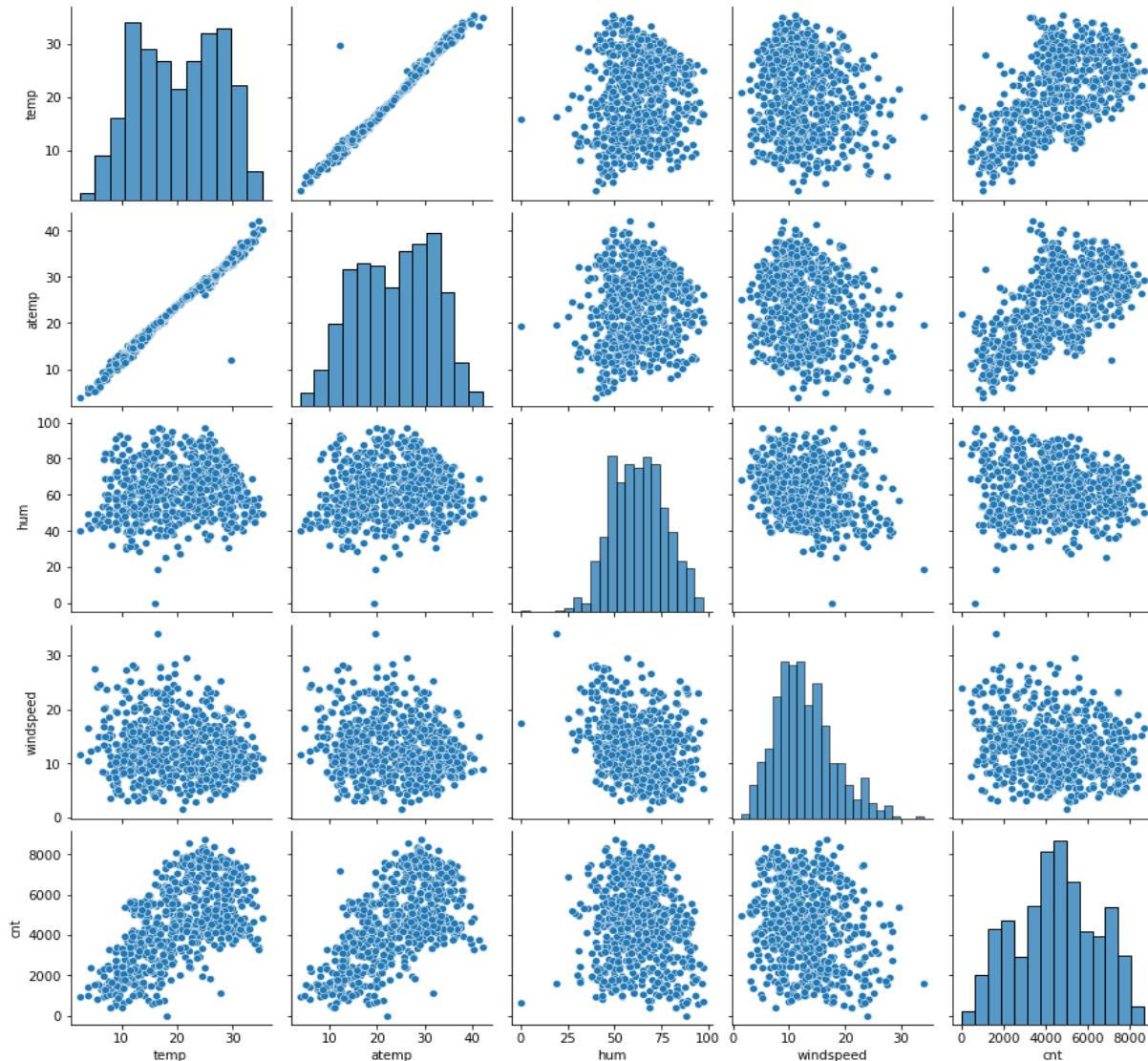
Ans:

This command is used to create (level -1) number of dummy variables Instead of creating same number of dummy variables as many as levels are there under a categorical variable . Thus the number of unnecessary column can be reduced.

Exmp: Suppose there are 3 levels named men, women and children under category variable. 10 for men & 01 for women is assigned then automatically 00 will denote rest one. Thus creating two dummies 0 & 1 for 3 levels easily can reduce one column for another dummy variable.

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

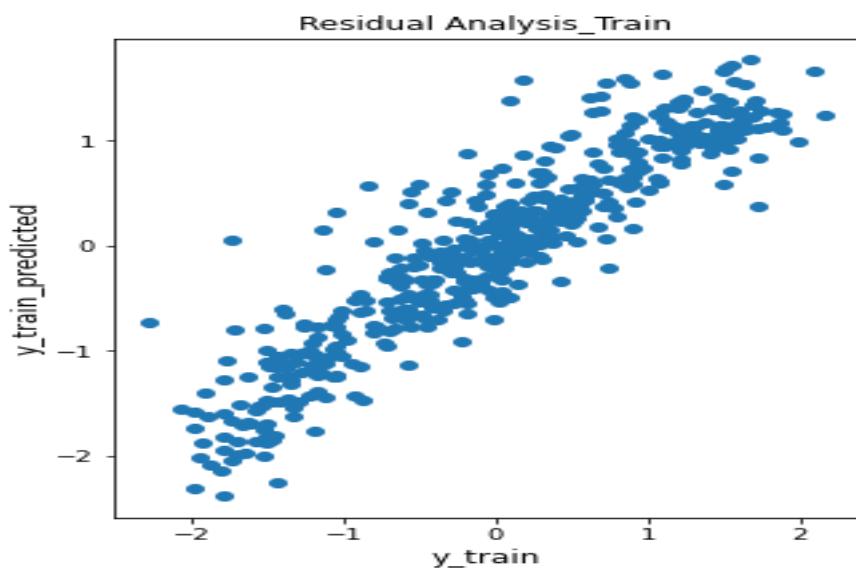
Ans: 'temp' and 'atemp' both variables have the highest correlation with target variable.



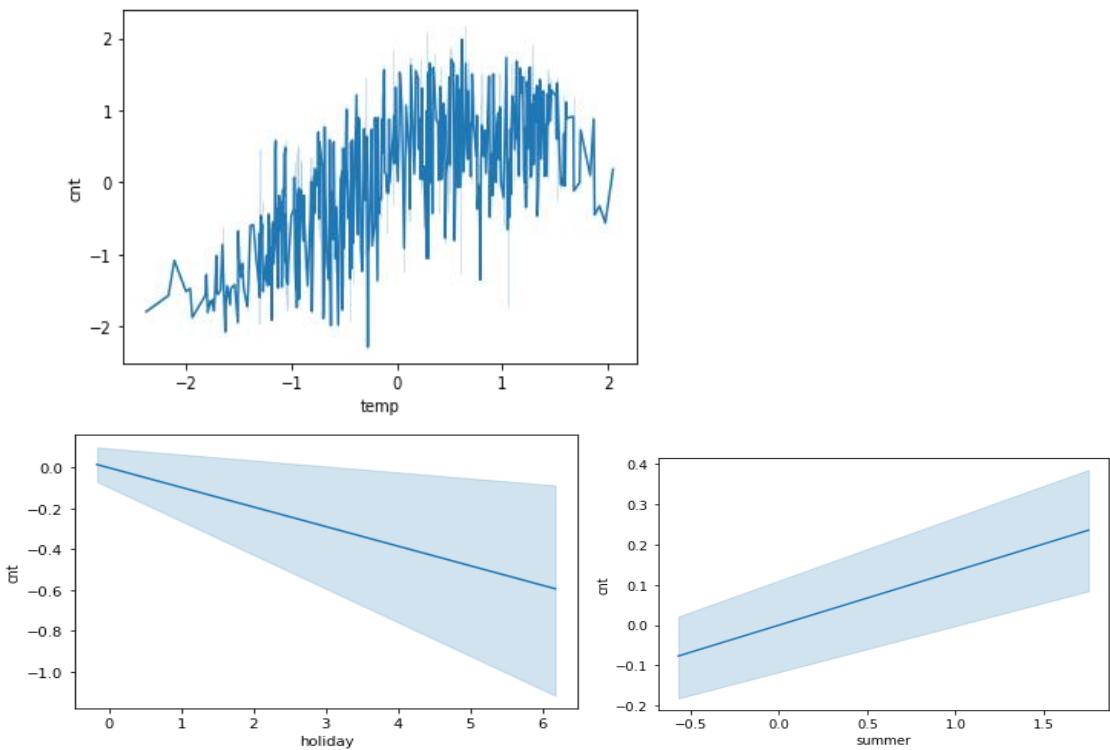
4. How did you validate the assumptions of Linear Regression after building the model on the training set? (3 marks)

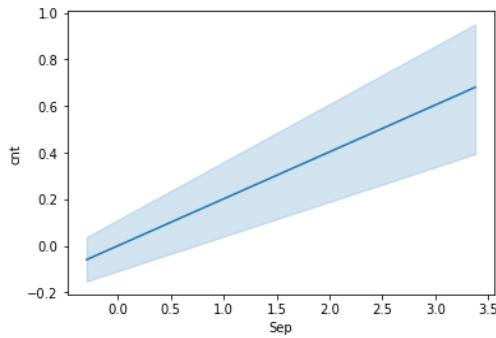
Ans:

- As Error points are randomly and almost uniformly distributed i.e. homoscedasticity has been found.

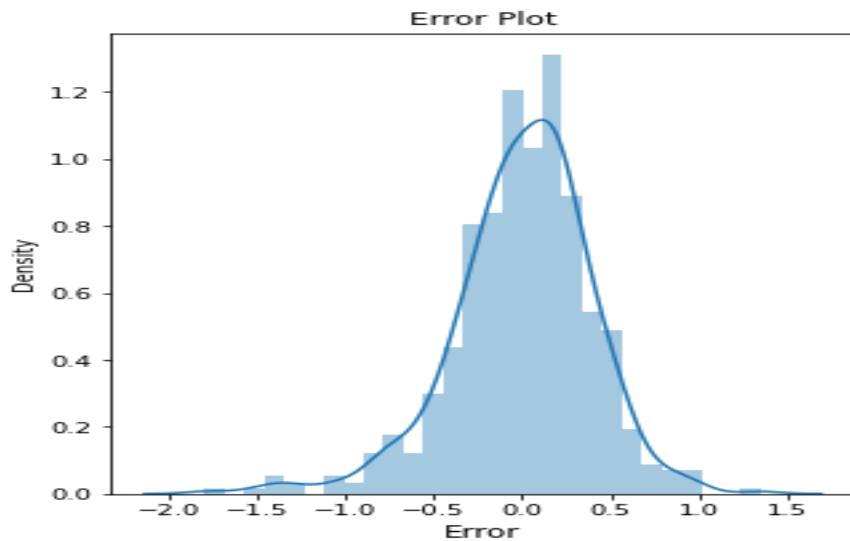


- Linear relationship found within dependent and independent variables.





3. Error terms are independent of each other.
4. Error terms are normally distributed.



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:

Temperature, Year and Season are the top 3 features contributing significantly towards explaining the demand of the shared bikes.

General Subjective Questions

1. Explain the linear regression algorithm in detail.

Ans:

Linear Regression:-

1. Process of estimating relationship between variables.
2. Explain change in dependent variable with change in the values of predictors.
3. Uses: forecasting and predictⁿ.
4. Shows correlatⁿ, not causatⁿ.
5. A form of parametric regression.
6. Interpolatⁿ not extrapolatⁿ.

↓
Model is used to predict the value of a dependent variable w.r.t. independent variables lying within the data range which already known.

↓
Model is used to predict the value of a dependent variable w.r.t. independent variable's value that lie outside the given range of the data, the model was built on.

A parametric model can be described using a finite number of parameters.

ex:- Linear Regression model build on n independent variables will have exactly n 'parameters' or coeff. by which entire model can be described.

Regression coeff:-

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

Arithmetic Mean of regressionⁿ is \geq correlatⁿ coeff.

Geometric mean of regression coeffs is = correlatⁿ coeff.

~~\bar{x} & \bar{y} = mean of given data~~

~~x_i & y_i = Observed data~~

Assumptions of Linear Regression :-

1. Linear Relationship between X & Y .
2. Error terms are normally distributed (not X, Y)
3. " " " independent of each other.
4. " " " have constant variance. (homoscedasticity)

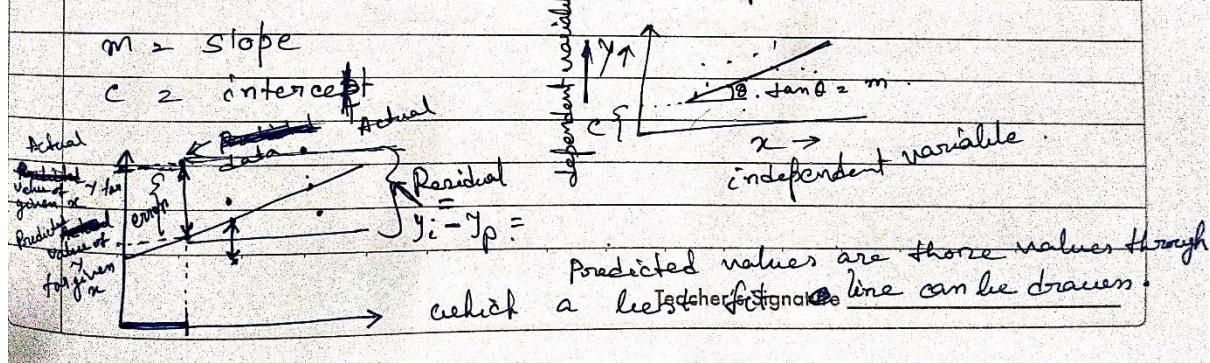
Linear Regression: → Supervised method.

→ Linear relationship between a dependent & one or more independent variables.

→ Used to find best fit line among the data.

→ " " " the coefficients of a line (m & c).
 $y = mx + c$. from a plot of actual data to find the values of predicted data. using $y_p = mx$

$$\text{predicted value } (y_p) = \text{Actual value of } x \text{ in } m + c \text{ independent variable}$$



Linear Regression! → Use Least square method.

$$r_i = y_i - y_p$$

Residual Sum of Squares (RSS) $\sum_{i=1}^n r_i^2$.

$$y_p = mx_i + c.$$

m & c derived from the plot of x_i and y_i .

$$RSS = (y_1 - c - mx_1)^2 + (y_2 - c - mx_2)^2 + \dots + (y_n - c - mx_n)^2$$

$$\boxed{RSS = \sum_{i=1}^n (y_i - c - mx_i)^2}$$

$\stackrel{\text{Actual}}{=} \sum_{\text{Error}}^{\text{Sum of Square of}}$

$$\Rightarrow \sum_{i=1}^n (y_i - y_{\text{predicted}})^2$$

2. Explain the Anscombe's quartet in detail.

Date _____

Expt. No. _____

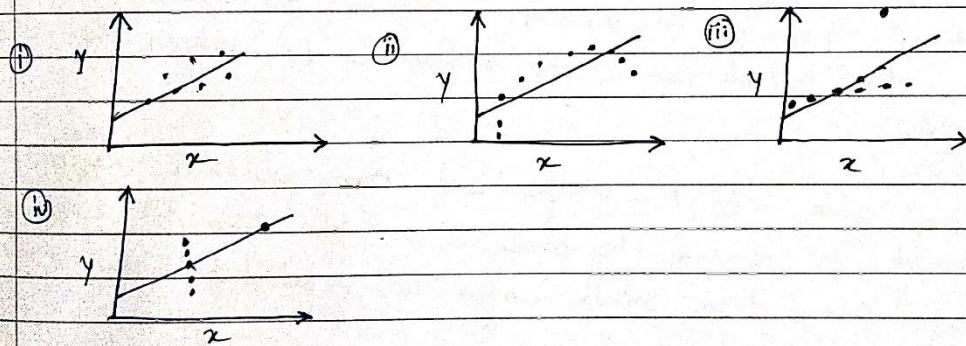
Page No. _____

Anscombe's quartet :- This consists of four data sets which have nearly identical simple descriptive statistics but very different distributions appeared when graphically plotted.

Therefore, Anscombe's quartet implies the importance of graphical analysis of data before the statistical analysis and model building applying any algorithm.

In this quartet, constructed by statistician Francis Anscombe, 11 (x, y) points are present for each dataset. Each dataset has almost identical statistic summary.

When the points are plotted graphical views are as follows:-



(i) The 1st scatter plot represents the linear relationship where y can be modelled as gaussian with mean linearly dependent on x .

(ii) The 2nd plot implies that the relationship between $x \& y$ is not linear so pearson correlatⁿ coeff. is not relevant. Something more efficient than general regression and its corresponding coeff. of correlatⁿ will be more appropriate.

(iii) In the 3rd case though the relatin between x & y is linear still a robust regression is needed to model it as there is outlier which has great influence to lower the correl. coeff.

(iv) In case of 4th graph we can see that the ^{one} high-lever point is enough to produce a high correlatⁿ coeff even though other data points show no relationship between them.

Thus the quarter shows how can the basic statistic properties make us fool based on the effect of outliers and other ~~and~~ inadequencies in data. So graphical analysis is must before applying any algorithm and building model.

3. What is Pearson's R?

Pearson's Coeff: -

Pearson Coeff. is a correlatⁿ coeff. used in statistic to determine the value of the linear relatⁿship between two variables. It is denoted by 'r'.

This rely on the hypothesis that the given data can be represented by a straight line. This coeff. is the rescaled version of regression coeff.

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

x_i = values of x variable in a sample

y_i = " " " "

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

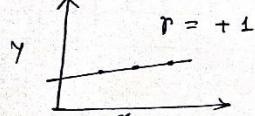
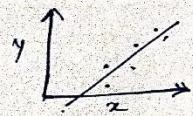
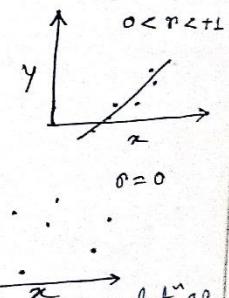
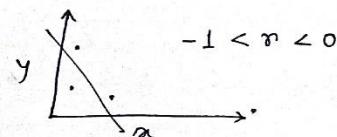
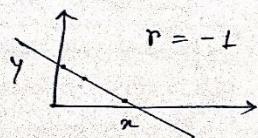
\bar{y} = " " " " Teacher's Signature — " —

Value of r varies in the range of -1 to $+1$.

If $r = -1$ then it means all data points lie on a line and the value of y decreases as the value of x increases or vice-versa.

If $r = +1$ then it means all data points lie on a line and the value of y increases as the value of x increases.

If $r = 0$ then there is no linear relationship between variables.



Hence this coeff. determine the strength of linear relatⁿ as well as directⁿ. The sign is determined by the regression slope (m).

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

Ans:

Scaling is the process of changing the range of data of a given dataset.

There are two methods of scaling:

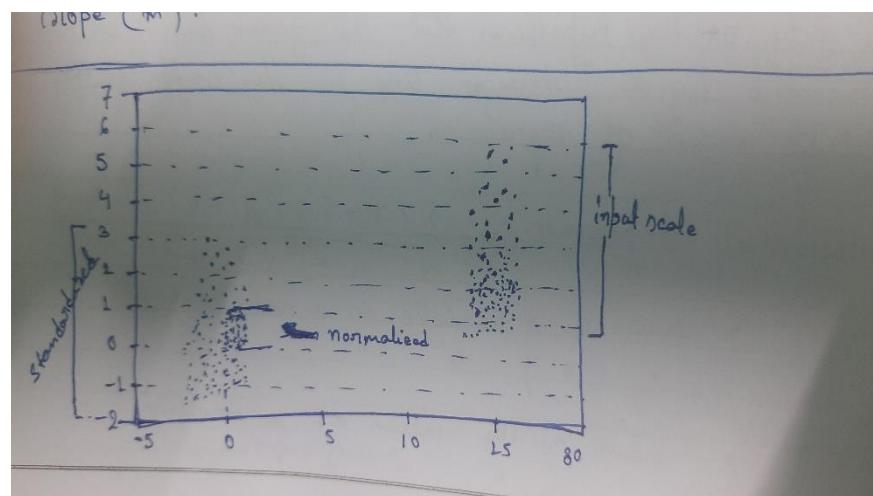
1. Normalization or Min-Max scaling
2. Standardization (mean-0, sigma-1)

It is extremely important to rescale the variables so that they have a comparable scale. If we don't have comparable scales, then some of the coefficients as obtained by fitting the regression model, might be very large or very small as compared to the other coefficients. This can affect model's performance. So it is advised to use standardization or normalization so that the units of the coefficients obtained are all on the same scale.

Due to rescaling optimization in the backend becomes faster also.

Standardization method doesn't compress data within a particular range as Normalization method does.

Standardization method is preferred over normalization in case of such data which have outliers



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

Ans:

Page No. _____

If there is a perfect relationship between two independent variables then we can get $R^2 = 1$. where R^2 = square of residual or correlatⁿ of coeff of determinatⁿ.

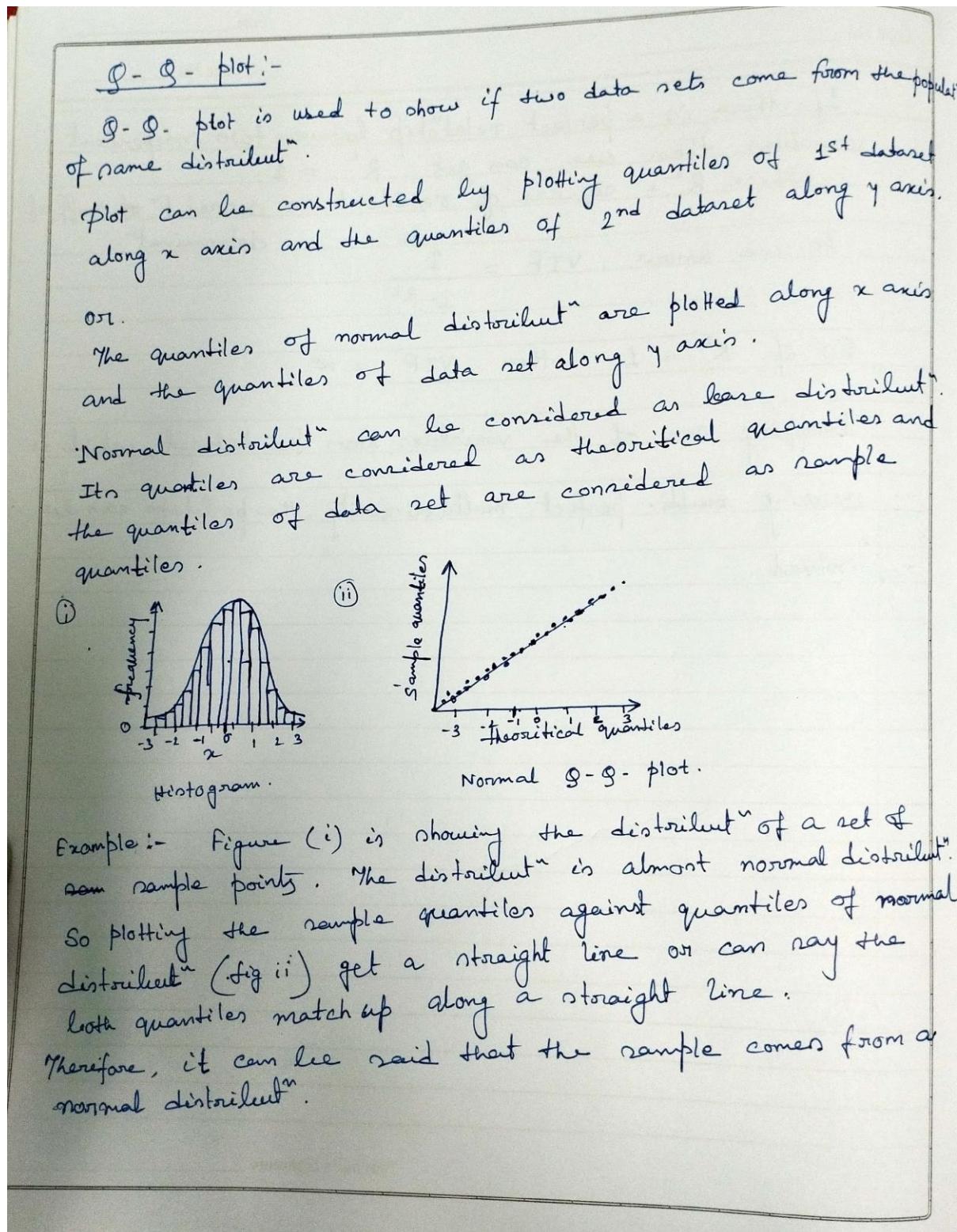
An we know $VIF = \frac{1}{1-R^2}$

So if $R^2 = 1$ then $VIF = \infty$

Dropping one of the variables, from the dataset, which is causing multi. perfect multicollinearity, the problem can be solved.

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

Ans:



Some advantages are there in using Q-Q plot like :-

1. The sample sizes do not need to be equal.
2. Many distributional aspects can be simultaneously tested, like,
 - a) shifts in locatⁿ.
 - b) shifts in scale.
 - c) changes in symmetry.
 - d) presence of outlier.

Interpretation:-

1. If all points of quantiles lie on or close to a straight line constructed at angle 45° , then the two data sets are coming from same distributⁿ.
2. If otherwise different distributⁿ is followed by the samples.

Use of importance in case of linear regression:-

- a) In case of linear regression, the train & test dataset can be examined by Q-Q-plot to determine if those are coming from populatⁿ of same distributⁿ.
- b) It helps to determine whether the residuals are normally distributed or not. If those are not aligne with the line with angle 45° then those are not following normal distributⁿ.

Teacher's Signature _____