# Predicting Rental Bikes

**Monesh soni**

**12 April 2020**

# Contents

## 1.1 Problem Statement

Predicting Bike rental count on daily based on the environmental and seasonal settings.

The aim of our project is predicting bike rental based on available data weather based or daily

registered and casual customer based on thus data we predict Bike rental count on daily based .

## 1.2 Data

Our task is to build model which count of rental bike based on weather and season daily based.

Our task is to build Regression model which will give the daily count of rental bikes based on weather and season Given below is a 10 observation sample of the data set that we are using to predict the count:

Table 1.1: Bike Rental Sample Data

| | instant | dteday | season | yr | mnth | holiday | weekday | workingday | weathersit | temp | atemp | hum | windspeed | casual | registered | cnt |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 1/1/2011 | 1 | 0 | 1 | 0 | 6 | 0 | 2 | 0.344167 | 0.363625 | 0.805833 | 0.160446 | 331 | 654 | 985 |
| 1 | 2 | 1/2/2011 | 1 | 0 | 1 | 0 | 0 | 0 | 2 | 0.363478 | 0.353739 | 0.696087 | 0.248539 | 131 | 670 | 801 |
| 2 | 3 | 1/3/2011 | 1 | 0 | 1 | 0 | 1 | 1 | 1 | 0.196364 | 0.189405 | 0.437273 | 0.248309 | 120 | 1229 | 1349 |
| 3 | 4 | 1/4/2011 | 1 | 0 | 1 | 0 | 2 | 1 | 1 | 0.200000 | 0.212122 | 0.590435 | 0.160296 | 108 | 1454 | 1562 |
| 4 | 5 | 1/5/2011 | 1 | 0 | 1 | 0 | 3 | 1 | 1 | 0.226957 | 0.229270 | 0.436957 | 0.186900 | 82 | 1518 | 1600 |
| 5 | 6 | 1/6/2011 | 1 | 0 | 1 | 0 | 4 | 1 | 1 | 0.204348 | 0.233209 | 0.518261 | 0.089565 | 88 | 1518 | 1606 |
| 6 | 7 | 1/7/2011 | 1 | 0 | 1 | 0 | 5 | 1 | 2 | 0.196522 | 0.208839 | 0.498696 | 0.168726 | 148 | 1362 | 1510 |
| 7 | 8 | 1/8/2011 | 1 | 0 | 1 | 0 | 6 | 0 | 2 | 0.165000 | 0.162254 | 0.535833 | 0.266804 | 68 | 891 | 959 |
| 8 | 9 | 1/9/2011 | 1 | 0 | 1 | 0 | 0 | 0 | 1 | 0.138333 | 0.116175 | 0.434167 | 0.361950 | 54 | 768 | 822 |
| 9 | 10 | 1/10/2011 | 1 | 0 | 1 | 0 | 1 | 1 | 1 | 0.150833 | 0.150888 | 0.482917 | 0.223267 | 41 | 1280 | 1321 |

As you can see in the table below we have the following 14 variables, using which we have to correctly predict daily based bike rental :

Table 1.2: Bike Rental Predictors variable

| s.no | Variables |
|---|---|
| 1 | Dteday |
| 2 | Season |
| 3 | Yr |
| 4 | Mnth |
| 5 | Holiday |
| 6 | Weekday |
| 7 | Workingday |
| 8 | Weathersit |
| 9 | Temp |
| 10 | Atemp |
| 11 | Hum |
| 12 | Windspeed |
| 13 | Casual |
| 14 | Registered |

Describe of dataset :

Describe returns the count, mean, standard deviation, minimum and maximum values and the quantiles of the Bike rental dataset.

Bike rental sample data Table 1.3 column(1-12)

| | instant | season | yr | mnth | holiday | weekday | workingday | weathersit | temp | atemp | hum | windspeed |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| count | 731.000000 | 731.000000 | 731.000000 | 731.000000 | 731.000000 | 731.000000 | 731.000000 | 731.000000 | 731.000000 | 731.000000 | 731.000000 | 731.000000 |
| mean | 366.000000 | 2.496580 | 0.500684 | 6.519836 | 0.028728 | 2.997264 | 0.683995 | 1.395349 | 0.495385 | 0.474354 | 0.627894 | 0.190486 |
| std | 211.165812 | 1.110807 | 0.500342 | 3.451913 | 0.167155 | 2.004787 | 0.465233 | 0.544894 | 0.183051 | 0.162961 | 0.142429 | 0.077498 |
| min | 1.000000 | 1.000000 | 0.000000 | 1.000000 | 0.000000 | 0.000000 | 0.000000 | 1.000000 | 0.059130 | 0.079070 | 0.000000 | 0.022392 |
| 25% | 183.500000 | 2.000000 | 0.000000 | 4.000000 | 0.000000 | 1.000000 | 0.000000 | 1.000000 | 0.337083 | 0.337842 | 0.520000 | 0.134950 |
| 50% | 366.000000 | 3.000000 | 1.000000 | 7.000000 | 0.000000 | 3.000000 | 1.000000 | 1.000000 | 0.498333 | 0.486733 | 0.626667 | 0.180975 |
| 75% | 548.500000 | 3.000000 | 1.000000 | 10.000000 | 0.000000 | 5.000000 | 1.000000 | 2.000000 | 0.655417 | 0.608602 | 0.730209 | 0.233214 |
| max | 731.000000 | 4.000000 | 1.000000 | 12.000000 | 1.000000 | 6.000000 | 1.000000 | 3.000000 | 0.861667 | 0.840896 | 0.972500 | 0.507463 |

Bike rental sample data Table 1.4 column(12-15)

| casual | registered | cnt |
|---|---|---|
| 731.000000 | 731.000000 | 731.000000 |
| 848.176471 | 3656.172367 | 4504.348837 |
| 686.622488 | 1560.256377 | 1937.211452 |
| 2.000000 | 20.000000 | 22.000000 |
| 315.500000 | 2497.000000 | 3152.000000 |
| 713.000000 | 3662.000000 | 4548.000000 |
| 1096.000000 | 4776.500000 | 5956.000000 |
| 3410.000000 | 6946.000000 | 8714.000000 |

# Chapter 2
# Methodology

## 2.1 Pre Processing

Any predictive modeling requires that we look at the data before we start modeling. However, in data mining terms looking at data refers to so much more than just looking. Looking at data refers to exploring the data, cleaning the data as well as visualizing the data through graphs and plots. This is often called as Exploratory Data Analysis. To start this process we will first try and look at all the probability distributions of the variables. Most analysis like regression, require the data to be normally distributed. We can visualize that in a glance by looking at the probability distributions or probability density functions of the variable.

### 2.1.1 Univariate Analysis

Univariate analysis is perhaps the simplest form of statistical analysis.The key fact is that only one variable is involved in univariate analysis.

In Figure 2.1 and 2.2 we have plotted the probability density functions numeric variables present in the data including target variable cnt.

i. Target variable cnt is normally distributed
ii. Independent variables like 'temp','atemp', and 'regestered' data is distributed normally.
iii. Independent variable 'casual' data is slightly skewed to the right so, there is chances of getting outliers.
iv. Other Independent variable 'hum' data is slightly skewed to the left , here data is already in normalize form so outliers are discarded.

Figure 2.1 Distribution of target variable (CNT) (python code in Appendix B)



Figure 2.2 showing distribution of dependent variables (python code in Appendix B)

**2.1.2 Bivariate   Analysis**
Bivariate  analysis  is  one  of  the  simplest  forms  of  quantitative  analysis.Bivariate  analysis  involve  two
variable. Bivariate analysis  helpful in testing simple hypotheses of association.
Let's we check now numeric variable temp or target variable cnt Bivariate relationshipBelow  ggpair
graph is showing clearly that relationship between  independent  variables 'temp' and  'atemp' are very
strong.

 1.The relationship between  'hum' , 'windspeed' with target variable 'cnt'


Figure 2.3 relationship between  numeric variables  (python  code in Appendix B)



Scatter plot is showing bivariate relationship between temp or cnt variable is high positive relationship

2.Relation between Numerical Variable 'hum' and target variable 'cnt'

 Figure 2.4 relationship between  numeric variables (python  code in Appendix B)

scatter plot is showing average relationship between hum variable or target variable cnt

**3. Relation between Numerical Variable 'atemp' and target variable 'cnt'**

Figure 2.5 relationship between numeric variables (python code in Appendix B)



Scatter plot showing highly positive relationship between "atemp" variable or target variable "cnt"

4. Bivariate relationship weekday is catogaries data or target variable cnt box plot to show distributions with respect to categories. box plot 'Weekdays' with 'cnt'

Figure 2.6 relationship between categories variables (python code in Appendix B)

**Boxplot shows most median value between 4000-5000**

5.Bivariate relationship : holiday is catogaries data or target variable cnt box plot to show distributions with respect to categories

box plot 'holiday' with 'cnt'

Figure 2.7 relationship between categories variables (python code in Appendix B)

## 2.2 Pre-processing technique

It is a data mining technique that transforms raw data into an understandable format. Raw data(real world data) is always incomplete and that data cannot be sent through a model. That would cause certain errors. That is why we need to preprocess data before sending through a model.

### 2.2.1 Missing Value Analysis

Missing values in data is a common phenomenon in real world problems. Knowing how to handle missing values effectively is a required step to reduce bias and to produce powerful models.

missing values, occur when no data value is stored for the variable in an observation. Missing data are a common occurrenceand can have a significant effect on the conclusions that can be drawn from the data.

Below table illustrate no missing value present in the Bike rental dataset.

 missing values Table

|  | Total | Percent |
|---|---|---|
| cnt | 0 | 0.0 |
| registered | 0 | 0.0 |
| casual | 0 | 0.0 |
| windspeed | 0 | 0.0 |
| hum | 0 | 0.0 |
| atemp | 0 | 0.0 |
| temp | 0 | 0.0 |
| weathersit | 0 | 0.0 |
| workingday | 0 | 0.0 |
| weekday | 0 | 0.0 |
| holiday | 0 | 0.0 |
| mnth | 0 | 0.0 |
| yr | 0 | 0.0 |
| season | 0 | 0.0 |
| dteday | 0 | 0.0 |
| instant | 0 | 0.0 |

## 2.2.2 Outlier    Analysis

The  Other steps  of Preprocessing Technique is Outliers analysis , an outlier is an observation point that is distant from other observations. Outliers in data can distort predictions and affect the accuracy, if you don't detect and handle them appropriately especially in regression models..

 As we are observed in fig 2.2  the data  is skewed so,  there is chance of outlier  in  independent variable 'casual' , one of the  best method to detect outliers is  Boxplot

Fig 2.4 shows   presence of Outliers in  variable 'casual'   and relationship between  'casual'  and 'cnt' before removing Outliers

Fig 2.5   shows  boxplot  of 'casual' after removing  outliers  and relationship between  'casual' and 'cnt' after removing  outliers

---

Boxplot :-  boxplot is a method for graphically depicting groups of numerical data through their quartiles. Box plots may also have lines extending vertically from the boxes (whiskers) indicating variability outside the upper and lower quartiles

Figure 2.4  'casual'  Baxoplot and  relation between 'cnt' and' casual' (R or python code in Appendix )

Below Boxplot in after removing outliers in casual variable

Figure 2.5 'casual' Boxplot   and relation between 'casual' and 'cnt' (R or python code in Appendix)



Since there is significant   difference  between  Pearson coefficient  correlation between  before  and after  outlier   detection  for  'casual' and 'cnt'  and losing  nearly  40 observation

## 2.2.3  Features  Selections

Feature selection or variable subset selection, is the process of selecting a subset of relevant.

Machine learning works on a simple rule – if you put garbage in, you will only get garbage to come out. By garbage here, I mean noise in data.

This becomes even more important when the number of features are very large. You need not use every feature at your disposal for creating an algorithm. You can assist your algorithm by feeding in only those features that are really important. I have myself witnessed feature subsets giving better results than complete set of feature for the same algorithm or – "Sometimes, less is better!".

Why should we require feature selection

- simplification of models to make them easier to interpret by researchers/users
- shorter training times.
- to avoid the curse of dimensionality,
- enhanced generalization by reducing overfitting (formally, reduction of variance)

We should consider the selection of feature for model  based on below criteria

i.        The relationship between two independent variable should  be  less and
ii.       The relationship between Independent and Target variables should be high.

Below fig 2.6  illustrates that relationship between  all numeric  variables using heatmap plot .

Figure 2.6  correlation plot of  variables (python code in Appendix B)



Color dark  red  indicates  there is strong  positive relationship and  if darkness is decreasing  indicates relation between  variables  are decreasing.

Color dark blue indicates there is strong negative relationship and if darkness is decreasing indicates relationship between variables are decreasing.

Above correlation plot is showing "temp" or "atemp" variable is highly positive correlated the are red color represented in above plot or "hum" variable highly negative relationship between target variable "cnt"

### 2.2.4 Dimensionality Reduction for variables

Above Fig 2.6 is showing

there is strong positive relationship between independent variables 'temp' and 'atemp'  so considering any one feature enough to predict the better Bike rental count .

And it is also showing there is almost no relationship between independent variable 'hum' and dependent variable 'cnt'.  so, 'hum' is not so important to predict Bike rental count .

So we drop thus two independent variable "atemp" or "hum" variable in our dataset using of drop() function .

## 2.2.5  Features  Scaling

The word "normalization" is used informally in statistics, and so the term normalized data can have multiple meanings. In most cases, when you normalize data you eliminate the units of measurement for data, enabling you to more easily compare data from different places. Some of the more common ways to normalize data include:

Transforming data using a z-score or t-score. This is usually called standardization. In the vast majority of cases, if a statistics textbook is talking about normalizing data, then this is the definition of "normalization" they are probably using.

Rescaling data to have values between 0 and 1. This is usually called feature scaling. One possible formula to achieve this is.

$$x_{new} = \frac{x - x_{min}}{x_{max} - x_{min}}$$

In rental dataset numeric variables like 'temp' , 'atemp' ,'hum' and ' windspeed' are in normalization form so , we have to Normalize two variables 'casual' and 'registered'

After normalize 'casual' and 'registered' variables look like in table below where all values between 0 and 1

Table  Normalization of 'casual' and 'registered

| casual | registered |
|---|---|
| 0.145833 | 0.091539 |
| 0.057181 | 0.093849 |
| 0.052305 | 0.174560 |
| 0.046986 | 0.207046 |
| 0.035461 | 0.216286 |

# Chapter 3
# Modelling

## 3.1 Model Selection

In earlier stage of analysis we have come to understand that few variables like 'temp' ,'casual, 'registered ' are going to play key role in model development for model development dependent variable may fall under below categories

i. Nominal
ii. Ordinal
iii. Interval
iv. Ratio

In our case dependent variable is interval so, the predictive analysis that we can perform is **Regression** Analysis

We will start our model building from Decision Tree .

### 3.1.1 Evaluating Regression Model

The main concept of looking at what is called **residuals** or difference between our predictions f(x[I,]) and actual outcomes y[i].

We are using three methods to evaluating performance of model

i. **MAPE** : The mean absolute percentage error, also known as mean absolute percentage deviation, is a measure of prediction accuracy of a forecasting method in statistics, for example in trend estimation, also used as a loss function for regression problems in machine learning. (Mean Absolute Percent Error) measures the size of the error in percentage terms. It is calculated as the average of the unsigned percentage error.

$$\left( \frac{1}{n} \sum \frac{|Actual - Forecast|}{|Actual|} \right) * 100$$

ii. **MSE: mean squared error** (**MSE**) or **mean squared deviation** (**MSD**) of an estimator (of a procedure for estimating an unobserved quantity) measures the average of the squares of the errors—that is, the average squared difference between the estimated values and the actual value. MSE is a risk function, corresponding to the expected value of the

squared error loss. The fact that MSE is almost always strictly positive (and not zero) is be because of randomness or because the estimator does not account for information that could produce a more accurate estimate.

$$\mathrm{MSE} = \frac{1}{n} \sum_{i-1}^{n} (Y_i - \hat{Y}_i)^2.$$

iii.     **RMSE :** (Root Mean Square Error) is a frequently used measure of the difference between values predicted by a model and the values actually observed from the environment that is being modelled. The **root-mean-square deviation** (**RMSD**) or **root-mean-square error** (**RMSE**) is a frequently used measure of the differences between values (sample or population values) predicted by a model or an estimator and the values observed. The RMSD represents the square root of the second sample moment of the differences between predicted values and observed values or the quadratic mean of these differences.

$$RMSE = \sqrt{\frac{\sum_{i=1}^{n} (X_{obs,i} - X_{model,i})^2}{n}}$$

## 3.2  Decision Tree

A tree has many analogies in real life, and turns out that it has influenced a wide area of **machine learning**, covering both **classification and regression**. In decision analysis, a decision tree can be used to visually and explicitly represent decisions and decision making. As the name goes, it uses a tree-like model of decisions.

Figure   3.2.1   Decision Tree Algorithm

Look at the above figure 3.2.1 here decision tree is using only two predictors variables to predict the model , which is not very impressive here the model is overfitted and biased towards only two predictors i.e 'casual' and 'registered' .

### 3.2.1 Evaluation of Decision Tree Model

Figure 3.2.2 Evaluation of Decision Tree using MAPE and RMSE

## Evaluation decison Tree regression model

In [103]:
```
1  #Define MAPE function
2  #Calculate Mean absolute percentage error (MAPE)
3  #we have create MAPE function to calculate
4  def MAPE(y_true, y_pred):
5      mape = np.mean(np.abs((y_true - y_pred) / y_true))*100
6      return mape
7
```

In [104]:
```
1  #Evaluate Decission tree regression  using of MAPE
2  #Mean absolute percentage error
3  MAPE(test_target_feature, DT_predictions)
4  #Here MAPE  is a measure of prediction accuracy of a forecasting method in statistics
5
```

Out[104]: 6.540085304302989

In [110]:
```
1  #Mean squared error
2  #we have import mean_squared_error to calculate MSE
3  from sklearn.metrics import mean_squared_error
4
5
```

In [111]:
```
1
2  #Mean squared error
3  mean_squared_error(test_target_feature,DT_predictions)
4  #Here MSE that is, the average squared difference between the estimated values and the actual value.
5
```

Out[111]: 37115.95652173913

```
1  #RMSE Root mean squared error
2  #math library available sqrt function
3  #sqrt function calculate square root
4  from math import sqrt
5  result=sqrt(mean_squared_error(test_target_feature,DT_predictions))
6  #here RMSE is a frequently used measure of the differences between values predicted
7  #by a model or an estimator and the values observed.
```

```
1  #Here RMSE
2  result
```

: 192.65501945638252

In Figure 3.2.3  Model Accuracy is   1- 6.54 = 0.934   which is nearly 93.4%  it is quite good  but  RMSE is 192.65  which is very high  so it's clearly stating that  our  Decision Tree Model is  Overfitted  and it working well for training data  but  won't predict good  for  new set of data. To overcome this overfit we have to tune the model using Random Forest Algorithm.

## 3.3 Random Forest

Random forest is a supervised learning algorithm. The "forest" it builds, is an ensemble of decision trees, usually trained with the "bagging" method. The general idea of the bagging method is that a combination of learning models increases the overall result.

Put simply: random forest builds multiple decision trees and merges them together to get a more accurate and stable prediction.

Random forest functions in below way

    i.       Draws a bootstrap sample from training data.
    ii.      For each sample grow a decision tree and at each node of the tree
          a.  Randomly draws a subset of mtry variable and p total of features that are available
          b.  Picks the best variable and best split from the subset of mtry variable
          c.  Continues until the tree is fully grown.

As we saw in section 3.2 Decision tree is overfitting and its accuracy MAPE and RMSE is also require to improve the performance of the regression model develop model using Random Forest Regressor we try to improve the model.

Figure 3.3.1 Random Forest Implementation

## Random Forest

```
1  #Implemantion of Random forest Regression Model
2  #Importing Random Forest Regressor Module from sklearn.ensemble pakage
3  from sklearn.ensemble import RandomForestRegressor
4  #train data set to store  object train_features_one indepandent or target variable
5  #We have to store variable in single object train_features_one independant variable
6  train_features=train[['season','yr','mnth','holiday','weekday','weathersit','temp','windspeed','casual','registered']].value
7  #Target variable
8  train_target_feature = train['cnt'].values
9
10 #test data set to store  object test_feature indepandent  or target variable
11 #We have to store variable in single object test_feature independant variable
12 test_feature=test[['season','yr','mnth','holiday','weekday','weathersit','temp','windspeed','casual','registered']].values
13 #Target variable
14 test_target_feature= test['cnt'].values
15 #Creating a Random Forest Regression model and fitting it to the training data
16 #A model comprised of many models is called an Ensemble model.
17 #For this model I've chosen 400 trees (n_estimator=400)
18 RF_model= RandomForestRegressor(n_estimators= 400, random_state=200).fit(train_features,train_target_feature)
```

```
1  #RF_model object
2  RF_model
```

```
: RandomForestRegressor(bootstrap=True, ccp_alpha=0.0, criterion='mse',
                        max_depth=None, max_features='auto', max_leaf_nodes=None,
                        max_samples=None, min_impurity_decrease=0.0,
                        min_impurity_split=None, min_samples_leaf=1,
                        min_samples_split=2, min_weight_fraction_leaf=0.0,
                        n_estimators=400, n_jobs=None, oob_score=False,
                        random_state=200, verbose=0, warm_start=False)
```

```
Random_state=200, verbose=0, warm_start=False)
```

```
1  #predict is a collection of objects, and the ordering of predictions is the same as the ordering of data passed in
2  RF_predictions= RF_model.predict(test_feature)
```

```
1  #Prediction of unseen data using or predict function
2  #prediction array value is
3  RF_predictions
```

Mtry :  Number of variables  to split at each node i.e. 7.

Node size :   size of each node  is 10

  Our  Random Forest model is looking  quite  good  where it utilized maximum  variables  to predict the count  values

**3.3.1 Evaluation of Random Forest**

 Figure  3.2.2 Random Forest Evaluation

## Evaluation Random Forest Model

```
1  #Using of MAPE fucntion we calculate Mean absolute percenatage
2  #Evaluate Random forest using  MAPE
3  MAPE(test_target_feature, RF_predictions)
4  #Mean absolute percentage error
```

2.6515610235581306

```
1  #Mean squared error
2  #we have import mean_squared_error to calculate MSE
3  from sklearn.metrics import mean_squared_error
4  mean_squared_error(test_target_feature,RF_predictions)
5  #Here MSE
6
```

12462.955450996378

```
1  #RMSE Root mean squared error
2  #math library available sqrt function
3  #sqrt function calculate square root
4  from math import sqrt
5  result=sqrt(mean_squared_error(test_target_feature,RF_predictions))
```

```
1  #RMSE
2  result
3  #Here Root mean squared error RMSE is given below
```

111.63760769111984

Fig 3.2.2  shows  Random Forest model performs  better  than Decision tree on both training and test data and well  also improve

(MAPE = 2.61) Accuracy =100-2.61: so model Accuracy =97.39

  and  RMSE decrease the  RMSE=(111) of the model  which is  quite impressive compression with Decision Tree model.

Using Linear Regression we will predict the Target variable 'cnt ' values and compare with Decision Tree model or Random Forest model.

## 3.4 Linear Regression

Multiple linear regression is the most common form of linear regression analysis. As a predictive analysis, the multiple linear regression is used to explain the relationship between one continuous dependent variable and two or more independent variables. The independent variables can be continuous or categorical.

**VIF ( Variance Inflation factor )** : the variance inflation factor is the quotient of the variance in a model with multiple terms by the variance of a model with one term alone. It quantifies the severity of multicollinearity in an ordinary least squares regression analysis. As Linear regression will work well if multicollinearity between the Independent variables are less.

Figure 3.4.1 Multi collinearity between Independent variables

```
The linear correlation coefficients ranges between:
min correlation ( mnth ~ yr ):  0.001342792
max correlation ( mnth ~ season ):  0.8198743

---------- VIFs of the remained variables --------
     Variables        VIF
1        season 3.887742
2            yr 2.807904
3          mnth 3.119420
4       holiday 1.122909
5       weekday 1.041991
6    workingday 3.257104
7     weathersit 1.355589
8          temp 2.547932
9     windspeed 1.115551
10       casual 3.805079
11   registered 6.785055
```

In the above figure it is showing there is strong correlation between two independent variable 'mnth' and 'season' so , it is enough to consider any one variable.

Figure 3.4.2  Multiple Linear Regression Model

## Linear Regression Model

```python
#We take here feature or data set
#train data set to store  object train_features_one independent or target variable
#We have to store variable in single object train_features_one independant variable
train_features = train[['season','yr','mnth','holiday','weekday','weathersit','temp','windspeed','casual','registered']].val
#Target variable
train_target_feature = train['cnt'].values

#test data set to store  object test_feature independant  or target variable
#We have to store variable in single object test_feature independant variable
test_feature = test[['season','yr','mnth','holiday','weekday','weathersit','temp','windspeed','casual','registered']].values
#Target variable
test_target_feature= test['cnt'].values
```

```python
#statsmodels is a Python module that provides classes and functions for the estimation of many different statistical models,
#as well as for conducting statistical tests, and statistical data exploration
import statsmodels.api as sm
#implement of Linear regression model
#OLS stands for Ordinary Least Squares
LR_model= sm.OLS(train_target_feature, train_features).fit()
```

```python
#summary of Linear Regression  model
#LR_model.summary()
#Using of summary function we get summary of Linear Regression model
LR_model
```

`<statsmodels.regression.linear_model.RegressionResultsWrapper at 0x2b6f6db8048>`

```python
#predicting Linear regression model using predict function
LR_predictions = LR_model.predict(test_feature)
```

### OLS Regression Results

| | | | |
|---|---|---|---|
| Dep. Variable: | y | R-squared (uncentered): | 1.000 |
| Model: | OLS | Adj. R-squared (uncentered): | 1.000 |
| Method: | Least Squares | F-statistic: | 6.604e+32 |
| Date: | Sat, 11 Apr 2020 | Prob (F-statistic): | 0.00 |
| Time: | 22:04:34 | Log-Likelihood: | 14223. |
| No. Observations: | 549 | AIC: | -2.843e+04 |
| Df Residuals: | 539 | BIC: | -2.838e+04 |
| Df Model: | 10 | | |
| Covariance Type: | nonrobust | | |

| | coef | std err | t | P>|t| | [0.025 | 0.975] |
|---|---|---|---|---|---|---|
| x1 | 1.847e-13 | 1.08e-13 | 1.703 | 0.089 | -2.84e-14 | 3.98e-13 |
| x2 | -1.307e-12 | 1.79e-13 | -7.318 | 0.000 | -1.66e-12 | -9.56e-13 |
| x3 | -9.948e-14 | 3.28e-14 | -3.034 | 0.003 | -1.64e-13 | -3.51e-14 |
| x4 | 1.705e-13 | 3.77e-13 | 0.453 | 0.651 | -5.69e-13 | 9.11e-13 |
| x5 | -1.279e-13 | 3.03e-14 | -4.214 | 0.000 | -1.88e-13 | -6.83e-14 |
| x6 | -8.527e-14 | 1.05e-13 | -0.814 | 0.416 | -2.91e-13 | 1.21e-13 |
| x7 | -5.002e-12 | 4.88e-13 | -10.246 | 0.000 | -5.96e-12 | -4.04e-12 |
| x8 | -1.364e-12 | 6.54e-13 | -2.088 | 0.037 | -2.65e-12 | -8.05e-14 |
| x9 | 1.0000 | 1.56e-16 | 6.39e+15 | 0.000 | 1.000 | 1.000 |
| x10 | 1.0000 | 7.03e-17 | 1.42e+16 | 0.000 | 1.000 | 1.000 |

| | | | |
|---|---|---|---|
| Omnibus: | 19.641 | Durbin-Watson: | 1.075 |
| Prob(Omnibus): | 0.000 | Jarque-Bera (JB): | 20.895 |
| Skew: | 0.475 | Prob(JB): | 2.90e-05 |
| Kurtosis: | 3.109 | Cond. No. | 4.63e+04 |

Here :

**Residual standard error**: 3.231e-12 on 576 degrees of freedom

**Multiple R-squared**:     1,        Adjusted R-squared:     1

 Here residual Standard error is quite less so  the distance between  predicted values f(x[I,]) and  actaual values f(x) are very less  so this model is predicted almost accurate values.

 And Multiple R-Square  value is 1 so, we can explain about 100 % of the data using our multiple linear regression model. This is very impressive in Linear Regression Model.

### 3.4.1 Evaluation of Linear regression Model

Figure  3.4.2  Evaluation of Regression Model

## Evaluation Linear Regression Model

```
1  #Evaluate Linear regression model   using of MAPE
2  #Mean absolute percentage error
3  MAPE(test_target_feature, LR_predictions)
4  #MAPE
```

5.591059667280465e-14

```
1  #Evaluate MSE
2  #Mean squared error
3  mean_squared_error(test_target_feature,LR_predictions)
4  #Here MSE =2.4407868961336697e-23
```

2.2227903388670907e-24

```
1  #Evaluate RMSE
2  #RMSE Root mean squared error
3  #math library available sqrt function
4  #sqrt function calculate square root
5  from math import sqrt
6  result=sqrt(mean_squared_error(test_target_feature,LR_predictions))
```

```
1  result
```

1.4909025249381968e-12

From above figure  it is clearly  showing  that  Model Accuracy is 99.9 % and  RMSE is nearly equal to 1.49.


## Model Selection

As we predicted counts  for Bike Rental using  three Models  Decision Tree, Random Forest and  Linear Regression  as MAPE is  high and RMSE is less for the Linear  regression  Model


**Conclusion**: - For the Bike Rental Data Linear Regression Model is   best model to predict the  count rental Bike daily based on environmental or seasonal .

**Appendix A -R Code**

# R code :

```r
#clear previously available list
rm(list = ls())
#Install required packages for implement of Bike renting prediction
#Load Libraries
#we have to store vector in x object
x = c("ggplot2", "corrgram", "DMwR", "caret", "randomForest", "unbalanced", "C50", "dummies", "e1071", "Information",
      "MASS", "rpart", "gbm", "ROSE", "sampling", "DataCombine", "dplyr","plyr","reshape","data.table","GGally","usdm")
#we have install packages using of command install.packages()
install.packages(x)
###############################################################################
#Require library to implement project
#library short description
# ggplot2 library is a system for declaratively creating graphics, based on The Grammar of Graphics
library("ggplot2")
#corrgram is Calculates correlation of variables and displays the results graphically
library("corrgram")
#Functions and data for "Data Mining with R"
library("DMwR")
#This package provides a number of functions for exploring the impact of different sources of uncertainties (e.g.positional uncertainty)
#on performance of species distribution models (SDMs).
library('usdm')
#The caret package (short for Classification And REgression Training) contains functions to streamline the |
#model training process for complex regression
library("caret")
#Library  random forest algorithm works by aggregating the predictions made by multiple decision trees of varying depth.
library("randomForest")
#The R package unbalanced implements some well-known techniques for unbalanced classification tasks and
#provides a racing strategy to adaptively select the best methods for a given dataset
library("unbalanced")
#The C50 package contains an interface to the C5.0 classification model
library("C50")
#Title Create dummy/indicator variables flexibly and efficiently.
library("dummies")
# e1071: Misc Functions of the Department of Statistics, Probability Theory Group (Formerly: E1071)
library("e1071")
```

```r
#The Information package is specifically designed to perform data exploration by producing
#easy-to-read tables and graphs
library("Information")
#MASS: Support Functions and Datasets for Venables and Ripley's MASS
library("MASS")
#rpart package keeps track of something called the complexity of a tree. The complexity measure is a combination of the size of
#a tree and the ability of the tree to separate the classes of the target variable.
library("rpart")
#gbm package in R to fit gradient boosting model.
library("gbm")
#Generation of synthetic data by Randomly Over Sampling Examples (ROSE)
library("ROSE")
#In sampling module sample takes a sample of the specified size from the elements of
#x using either with or without replacement.
library("sampling")
#Using of Data combine module A function for filling in missing values of a variable from one
#data frame with the values from another variable
library("DataCombine")
# Dplyr aims to provide a function for each basic verb of data manipulation.
library("dplyr")
#plyr is an R package that makes it simple to split data apart, do stuff to it, and mash it back together.
#This is a common data-manipulation step.
library("plyr")
#Different functions require different formats, and so the need to reshape a dataset.
library("reshape")
#Data. table is an extension of data. frame package in R. It is widely used for fast aggregation of large datasets,
#low latency add/update/remove of columns, quicker ordered joins, and a fast file reader.
library("data.table")
#'GGally' module extends 'ggplot2' by adding several function or plot
library("GGally")
###############################################################################
```

```r
#get working directory
#using of get function we check current working diroctory
getwd()
#Set working directory
#Using of set function we set our working diroctory
setwd("C:/Users/ASUS/Desktop")

#Load Day dataset
Bike=read.csv("Day.csv",header =TRUE)
#Using of head function we get  five observation in our  Bike dataset
head(Bike)
#In Bike dataset available 16 variable
#In this variable 15 are indepandent variable or cnt is target variable
#Using of tail fucntion we get last five observation in our Bike dataset
tail(Bike)
#Measures of centrality are mean, median. Measures of spread are
#standard deviation, quantiles, min and max, range, interquartile range.
#Using of summary function we get  description of Bike object
summary(Bike)
#str() It provides great information about the structure of some object.
#str() function is provide great information of Object
str(Bike)
##############################################################################
#Univariate Analysis#

#Univariate analysis is the simplest form of analyzing data. "Uni" means "one",
#so in other words your data has only one variable. It doesn't deal with causes
#or relationships (unlike regression) and it's major purpose is to describe;
#it takes data, summarizes that data and finds patterns in the data

#glimpse function is helpful for a first inspection of the data set at hand.
glimpse(Bike)




#Univariate analysis in Numeric variable

#Analyze the distribution of  target variable "cnt"
hist(Bike$cnt,right = FALSE,col = "blue")
#Above histogram graphical cnt variable show normal distribution

#Analyze the distribution of indepandent variable "casual"
hist(Bike$casual,right = F,col = "blue")
#Aboove histogram graph showing casual variable is not normaly distributed and
#chances of outlier is high in this variable

#Analyze the distribution of indepandent variable "temp"
hist(Bike$temp,right = F,col = "blue")
#Aboove histogram graph showing temp variable is normaly distributed

#Analyze the distribution of indepandent variable "registered"
hist(Bike$registered,right = F,col = "blue")
#Aboove histogram graph showing ragistered variable is normaly distributed

#Analyze the distribution of indepandent variable "hum"
hist(Bike$hum,right = F,col = "blue")
#Aboove histogram graph showing hum variable is normaly distributed

##############################################################################
#Categorical variable Univariaite analysis

# Visualize categorical variable "mnth" with target variable "cnt"
#Change x varable factor catogroy using of as.factor function
#using of fill function we fill color in our plot
ggplot(Bike, aes(x=as.factor(mnth), y=cnt),fill="blue") +
  stat_summary(fun="mean", geom="bar",color="red")
#Using of aes() aesthatic function we fill color visualize plot just like (boxplot,histogram)
#These visual caracteristics are known as aesthetics (or aes) and include:
#color and fill. points shape.
```

```
ggplot(Bike)+
  geom_histogram(aes(x=cnt,y=..density..),
                 fill= "red")+
  geom_density(aes(x=cnt,y=..density..))

# Visualize categorical variable 'holiday'

ggplot(Bike) +
  geom_bar(aes(x=holiday),fill="red")
#Number of count is very high in holiday time renting bike
# it is showing that  the mostly cycle rentals are happening  on holidays

# Visualize categorical variable 'weekday'

ggplot(Bike) +
  geom_bar(aes(x=weekday),fill="red")
#Above plot is showing number of counts is same in all weekdays

# Visualize categorical variable 'weathersit'
ggplot(Bike) +
  geom_bar(aes(x=weathersit),fill="blue")
#Number of count is very high in whether is "clear"
#Number of count is average if whether "Few clouds"
#Number of count is less if whwther is  "Partly cloudy, Partly cloudy"

##########################################################################
#Bivariate Analysis

#Bivariate relationship in numeric variable
#Firt we check Bivariate relationship "temp" and "hum" variable
#Using of geom_pointThe point geom is used to create scatterplots.
#geom_smooth understands the following aesthetics (required aesthetics are in bold):
#x,y,alpha,colour,fill,linetype,size,weight

ggplot(Bike, aes(x= temp,y=hum)) +
  geom_point()+
  geom_smooth()
#It is showing  graph  humadity increasing rapidly till 0.65 after after the
#have decressing gredualy


#we check Bivariate relationship "temp" and "atemp" variable
ggplot(Bike, aes(x= temp,y=atemp)) +
  geom_point()+
  geom_smooth()
#It is showing graph "temp" or "atemp" varaible is strongly positive bivariate relationship
#temp or atemp variable is highly positive correlated

#we check Bivariate relationship "temp" and "windspeed" variable
ggplot(Bike, aes(x= temp,y=windspeed)) +
  geom_point()+
  geom_smooth()
#It is showing graph "temp" or "windspeed" varaible is less negative correlation
#tempraure is increases thus time windspeed is decreass

################################################################
#Checking relationship between all numeric variable
#We check  thus numeric relationship using of pair plot
ggpairs(Bike[,c('atemp','temp','hum','windspeed','cnt')])
#Above plot is showing positive correlation atemp or target variable cnt
#plot is showing positive correlation temp or target variable cnt
#plot is showing negative correlation hum or target variable cnt
#plot is showing positive correlation windspeed or target variable cnt

################################################################
```

```
#Catogirical Bivariate Relationship

#check relationship between  season and holiday
#table() performs categorical tabulation of data with the variable and its frequency.
rel_session_holiday= table(Bike$season,Bike$holiday)
#In barplot function We can supply a vector or matrix to this function.
barplot(rel_session_holiday)
#If a day is neither weekend nor holiday is 1, otherwise is 0.
#Here chances of renting working day is less or chances or renting
# holiday is very high

#check relationship between  season and weekand
rel_season_weekday=table(Bike$season,Bike$weekday)
barplot(rel_season_weekday)
#weekday bike renting chancs all day are same

#Check relationship between season and weathersit
rel_season_weathersit=table(Bike$season,Bike$weathersit)
rel_season_weathersit
barplot(rel_season_weathersit)
#weathersit is three catigory
# Clear, Few clouds, Partly cloudy
#Bike rental is very high in clear wether
#Bike rental is normal in Few cloud season wether
#Bike rental is less in partly cloudy wether


########################################################################


#Missing Value Analysis
# missing values, occur when no data value is stored for the variable in an observation
#WE find missing value in using of is.na function
#We calculate all missing value using of sum function
```

```
missing_value = data.frame(apply(Bike,2,function(x){sum(is.na(x))}))
missing_value$columns=row.names(missing_value)
#Object
missing_value
names(missing_value)[1] =  "Missing_value_percentage"
#Using of names function we define column name
missing_value$Missing_value_percentage = (missing_value$Missing_value_percentage/nrow(Bike)) * 100
#Above formula get percentage of missing value in our variable
missing_value = missing_value[order(-missing_value$Missing_value_percentage),]
row.names(missing_value) = NULL
missing_value = missing_value[,c(2,1)]

missing_value
 #In our Bike dataset is not available missing value

#################################################################

#Outlier Analysis
#we perform outlier Analysis in numeric variable
#Allready we done univariate analysis
#The chances of getting outlier more in "casual","ragistered"or cnt variable
#Outlier analysis in registered variable
boxplot(Bike$registered,col = "yellow")
#plot  showing there is not available outlier value

#Outlier Analysis in "casual" variable
boxplot(Bike$casual,col = "yellow")
#plot showing outlier available in casual variable

#Outlier Analysis in "cnt" variable
boxplot(Bike$cnt,col = "yellow")
#plot showing outlier not available in cnt variable

##########################################################################
```

```r
#Handling outliers
#Remove outlier using of Boxplot method
df=Bike
#Bike = df
#We remove outliers in "casual" variable
value=Bike$casual[Bike$casual%in% boxplot.stats(Bike$casual)$out]
Bike=Bike[which(!Bike$casual %in% value),]

#Now our Casual variable not available outliers
boxplot(Bike$casual,col = "red")
#Boxplot showing not available outliers in casual variable
#Before outliers number of bservation is 731 but after removing outliers observation is 687

############################################################################
#Feature selection
#Correlation Analysid
#we are perform correlation analysis using of corrgram library
corrgram(Bike[,c("temp","atemp","hum","windspeed","registered" ,'cnt')], order = F,
         upper.panel=panel.pie, text.panel=panel.txt, main = "Correlation Plot")
#Blue color show highly positive correlationship between variable
#Correlation plot is showing temp and atemp variable is highly correlated
#"hum" variable or target variable "cnt" not available relationship betweent them

############################################################################


#Reduction of variable the are highly correlated or not related
#we know that "temp" or "atemp" variable is highly correlated "hum" variable is not related with target variable "cnt"
#We are remove here "atemp" or "hum" variable

Bike_features = subset(Bike,select=-c(atemp,hum))
View(Bike_features)
#our Bike dataset now available 687 observation or 14 variable
```

```r
#Normalisation
#OUR DATASET  most numeric variable is normalization form
#Normalized variable is "temp","atemp","hum","windspeed"
#so we perform normalisation "casual" or "registered" variable

cnames = c("casual","registered")
for(i in cnames){
  print(i)
  Bike_features[,i] = (Bike_features[,i] - min(Bike_features[,i]))/
    (max(Bike_features[,i] - min(Bike_features[,i])))
}
Bike$casual
Bike$registered

############################################################################
#Model development
#Bike dataset available variable we store columns name object
columns=c("season" ,"yr" ,"mnth","holiday","weekday","workingday","weathersit","temp","windspeed","casual","registered","cnt")

#Divide data into train and test using stratified sampling method
#set.seed function in R is used to reproduce results i.e. it produces the same sample again and again. When we generate randoms numbers without se
#function it will produce different samples at different time of execution.
set.seed(1234)

#we split data intoo traing and test dataset
#train dataset available 80% of data
#test dataset available 20% of data
#we are using to perform operation createDatPartion() funtion the are available in caret() library
train.index = createDataPartition(Bike_features$cnt, p = .80, list = FALSE)
train = Bike_features[ train.index,]
test  = Bike_features[-train.index,]
train_feature = train[,c("season" ,"yr" ,"mnth","holiday","weekday","workingday","weathersit","temp","windspeed","casual","registered","cnt")]
```

```r
338  #train dataset
339  train_feature
340
341  test_features = test[,c("season" ,"yr" ,"mnth","holiday","weekday","workingday","weathersit","temp","windspeed","casual","registered","cnt")]
342  #test dataset
343  test_features
344
345  ##############################################################################
346
347  #Decision Tree for regression
348  #our dataset target variable is based on continuos prediction so we select Regresssion Algorithm
349  #rpart of regression
350  #rpart keeps track of something called the complexity of a tree. The complexity measure is a combination of the size of
351  #a tree and the ability of the tree to separate the classes of the target variable
352  fit = rpart(cnt ~ ., data = train_feature, method = "anova")
353
354
355  #Predict for new test cases
356  DT_predictions = predict(fit, test_features[,-12])
357
358  print(fit)
359
360  DT_predictions
361  plot(DT_predictions,col="blue")
362  #  plotting decision tree
363  #using of pair function The function 'par' can be  used  to  manipulate
364  #the  current plot( ) . pair function allow multiple plot in on one figure
365  #cex: for legends and other independent functions
366  #outside  of  the 'plot'  family
367  par(cex= 0.8)
368  #plot fit object
369  plot(fit)
370  #print text in regression decision tree
371  text(fit)


374
375  #Decision Tree Model Evaluation
376
377  #Root Mean Squared Error
378  #(RMSE): RMSE is a quadratic scoring rule that also measures the average magnitude of the error.
379  #Evaluate  Model using RMSE
380  #make RMSE function first the evaluate RMSE formula
381  RMSE <- function(y_test,y_predict) {
382
383      difference = y_test - y_predict
384      root_mean_square = sqrt(mean(difference^2))
385      return(root_mean_square)
386
387  }
388  RMSE(test_features[,12], DT_predictions)
389
390  #RMSE= 453.6728
391
392  ###############################################################
393
394  #MAPE Mean Absolute Percentage Error
395  #MAPE:is a measure of prediction accuracy of a forecasting method in statistics
396  #First we create MAPE function
397  #calculate MAPE
398  MAPE = function(y, y_true){
399      mean(abs((y - y_true)/y))
400      }
401
402  MAPE(test_features[,12], DT_predictions)
403  #error rate = 0.1037541
404  #Accuracy of model is 90.63
405
406  ##############################################################################
```

```r
408  #train or test dataset and selection of variable
409  train = Bike_features[ train.index,]
410  test  = Bike_features[-train.index,]
411  train_feature = train[,c("season" ,"yr" ,"mnth","holiday","weekday","workingday","weathersit","temp","windspeed","casual","registered","cnt")]
412  #train dataset
413  train_feature
414
415  test_features = test[,c("season" ,"yr" ,"mnth","holiday","weekday","workingday","weathersit","temp","windspeed","casual","registered","cnt")]
416  #test dataset
417  test_features
418
419
420  ################################################################################
421
422  #Random Forest Model
423  #In the random forest approach, a large number of decision trees are created.
424  #Every observation is fed into every decision tree.
425  RF_model=randomForest(cnt~.,data=train_feature,ntree=100)
426  RF_model
427  #Random forest model for train_feature dataset
428  plot(RF_model,col="red")
429  #plot of random forest model
430  #prediction of test case
431  RF_predictions = predict(RF_model, test_features[,-12])
432  RF_predictions
433  plot(RF_predictions,col="red")
434  #predicted value in test_feature dataset
435
436  ################################################################################
437
438  #Evaluation of Random Forest algorithm
439
440  #Root Mean Squared Error
441  #RMSE
442  ⌐
```

```r
442  RMSE(test_features[,12], RF_predictions)
443  #RMSE=207.7042
444  #MAPE Mean Absolute Percentage Error
445  MAPE(test_features[,12], RF_predictions)
446  #Error Rate :0.04592139
447  #Accuracy Rate :96.5
448
449  ####################################################################################
450
451  #train or test dataset and selection of variable
452  train = Bike_features[ train.index,]
453  test  = Bike_features[-train.index,]
454  train_feature = train[,c("season" ,"yr" ,"mnth","holiday","weekday","workingday","weathersit","temp","windspeed","casual","registered","cnt")]
455  #train dataset
456  train_feature
457
458  test_features = test[,c("season" ,"yr" ,"mnth","holiday","weekday","workingday","weathersit","temp","windspeed","casual","registered","cnt")]
459  #test dataset
460  test_features
461
462
463  #Linear Regression Model
464
465  #Variance Inflation Factor
466  #Calculates the variation inflation factors of all predictors in regression model
467  vif(train_feature[,-12])
468  vifcor(train_feature[,-12], th = 0.9)
469  # Correleation between two variables is 'season' and 'mnth' is 0.82 so,
470  #removing one variable from the model
471
472  #column variable name train dataset
473  train_feature = train[,c("yr" ,"mnth","holiday","weekday","workingday","weathersit","temp","windspeed","casual","registered","cnt")]
474  #Column variable name test dataset
475  test_features = test[,c("yr" ,"mnth","holiday","weekday","workingday","weathersit","temp","windspeed","casual","registered","cnt")]
```

```
476  # the lm(), or "linear model," function can be used to create a simple regression model
477  LR_model = lm(cnt ~., data = train_feature)
478  print(LR_model)
479
480  #Predictions of Linear Regression Model
481  LR_predictions = predict(LR_model, test_features[,-11])
482  print(LR_predictions)
483  plot(LR_predictions,col="red")
484
485▾ #########################################################################
486
487  #Evaluation of Linear Regression Model
488  ##Root Mean Squared Error
489  #RMSE
490  RMSE(test_features[,11], LR_predictions)
491  #RMSE 1.591802e-12
492
493  ##MAPE Mean Absolute Percentage Error
494  #MAPE
495  MAPE(test_features[,11], LR_predictions)
496  #Accuracy :99.9 +Accuracy
497  #Error rate : 1.065695e-15
498▾ #########################################################################
499
500  #Conclusion for Bike dataset Linear Regression model is a best model and give
501  #approximately 99.9 percent accuracy for Bike rental count on daily based
```

# Appendix B -Python Code

# Python code:

```python
1  #importing libraries
2  #The OS module in Python provides a way of using operating system dependent  functionality.
3  import os
4  #Pandas is a module use for data manipulation and analysis.
5  import pandas as pd
6  #Numpy library is useful linear algebra, Fourier transform,array,matrix and random number capabilities
7  import numpy as np  |
8  #Matplotlib is a comprehensive library for creating static, animated, and interactive visualizations in Python.
9  import matplotlib.pyplot as plt
10 from  matplotlib import pyplot
11 #SciPy contains modules for optimization, linear algebra, integration, interpolation, special functions
12 from scipy.stats import chi2_contingency
13 #Seaborn is a Python data visualization library based on matplotlib
14 import seaborn as sns
15 #random module depend on a pseudo-random number generator function random()
16 from random import randrange, uniform
17 #The datetime module supplies classes for manipulating dates and times
18 import datetime
19 #Scikit-learn or sklearn modual provides a range of supervised and unsupervised learning algorithms -
20 #via a consistent interface in Python
21 from sklearn.ensemble import RandomForestRegressor
22 from sklearn.model_selection import train_test_split
23 from sklearn.tree import DecisionTreeRegressor
24  #statsmodels is a Python module that provides classes and functions for the estimation of many different statistical models
25 #as well as for conducting statistical tests, and statistical data exploration
26 import statsmodels.api as sm
27
```

```python
1  #change working diroctory using chdir function available in os library ch=change ,dir=diroctory
2  os.chdir(r"C:\Users\ASUS\Desktop")
3
```

```python
1  #Loading dataset
2  #we have loading dataset using of pandas library or read csv file
3  Bike=pd.read_csv("day.csv")
```

```python
1  #top 10 observation we get using of head function Bike is object here
2  #head()"function of pandas library which returns first five observations of the data set.
3
4  Bike.head(10)
5
```

```python
1  #".tail()" returns last five observations of the data set.
2
3
4  Bike.tail(5)
```

```python
1  #Understanding of dataset
2  # ".shape" return number of observation row or column
3  Bike.shape
4  #Bike data set contain 731 row or 16 column
5
6  #The describe() function in pandas is very handy in getting various summary statistics.
7  #This function returns the count, mean, standard deviation, minimum and maximum values and the quantiles of the data.
8  Bike.describe()
```

```python
1  #data set contain integer ,float or variable dteday object(catogrical) type
2  #using of info() funtion we get information of our dataset
3  Bike.info()
4
```

```python
1  #all variable  name in dataset
2  col_name=Bike.columns
3  pd.DataFrame(col_name)
```

# Univariate Analysis

```
1  #Univariate analysis is perhaps the simplest form of statistical analysis.
2  #The key fact is that only one variable is involved in univariate analysis.
```

```
1
2  #Target variable  analysis
3  #descriptive statistics summary
4  #describe() function in pandas is very handy in getting various summary statistics
5  #cnt variable is our target variable
6  Bike['cnt'].describe()
7
```

```
1  #Check whether target variable is normal or not
2  #most value cnt variable is 0.00005 to 0.00015
3  #Flexibly plot a univariate distribution of observations.
4  sns.distplot(Bike['cnt'],color="m")
5  #Below plot is showing cnt variable is normaly distributed
```

```
1  #check atemp variable normal or not their available outlier value or not or range in numeric value using of seborn library
2  #distplot function plot to atemp  variable
3  sns.distplot(Bike['atemp'],color="m")
4
5
```

```
1  #check windspeed variable normal or not their available outlier value or not or range in numeric value using of seborn libra
2  #distplot function plot to windspeed  variable
3
4  sns.distplot(Bike['windspeed'],color="m")
```

```
1  #skewnwss is Assymetry distributation
2  #Here skewness  value is -0.047353 their is negative skewness
3  #Here Kurtosis value is -0.811922 if value is less than 0 than their is playtocartic
4  #Here kurtosis is very less so target variable is playticartic distribitation
5  print("Skewness: %f" % Bike['cnt'].skew())
6  print("Kurtosis: %f" % Bike['cnt'].kurt())
```

# Bivariate Relationship

```
1   #Bivariate analysis is one of the simplest forms of quantitative analysis.
2   #Bivariate analysis involve two variable.
3   # Bivariate analysis  helpful in testing simple hypotheses of association.
4   #Let's we check now numeric variable temp or target variable cnt Bivariate relationship
5   #Using of value_counts() function returns object containing counts of unique values.
6   Bike['temp'].value_counts()
7   #Now we drow scatter plot in betweent temp or cnt variable
8   variable = 'temp'
9   #using of concati function in pandas library
10  #Concatenate pandas objects along a particular axis with optional set logic along the other axes
11  Bi_ralationship= pd.concat([Bike['cnt'], Bike[variable]], axis=1)
12  #Using of matplotlib library scatter(x,y) creates a scatter plot with circles at the locations specified by
13  #the vectors x and y
14  Bi_ralationship.plot.scatter(x=variable, y='cnt', ylim=(0,10000),color="red")
15  #scatter plot is showing positive relationship between temp or cnt variable
16
17
```

```
1   #relation between Numerical Variable 'hum' and target variable 'cnt'
2   Bike['hum'].value_counts()
3   #Now we draw scatter plot between 'hum' and 'cnt' variables
4   variable = 'hum'
5   #Concatenate pandas objects along a particular axis with optional set logic along the other axes
6   Bi_ralationship= pd.concat([Bike['cnt'], Bike[variable]], axis=1)
7   #Using of matplotlib library scatter(x,y) creates a scatter plot with circles at the locations specified by
8   #the vectors x and y
9   Bi_ralationship.plot.scatter(x=variable, y='cnt', ylim=(0,10000),color="red")
10  #scatter plot is showing average relationship between hum variable or target variable cnt
11
12
```

```
1   #relation between Numerical Variable 'windspeed' and target variable 'cnt'
2   Bike['windspeed'].value_counts()
3   #Now we draw scatter plot between 'atemp' and 'cnt' variables
4   variable = 'windspeed'
5   #Concatenate pandas objects along a particular axis with optional set logic along the other axes
6   Bi_ralationship= pd.concat([Bike['cnt'], Bike[variable]], axis=1)
7   #Using of matplotlib library scatter(x,y) creates a scatter plot with circles at the locations specified by the vectors x an
8   Bi_ralationship.plot.scatter(x=variable, y='cnt', ylim=(0,10000))
9   #scatter plot is showing negative relationship between windspeed and target variable cnt
```

```
1   #Bivariate relationship  weekday is catogaries  data or target variable cnt
2   #box plot to show distributions with respect to categories
3   #box plot 'Weekdays' with 'cnt'
4   variable_weekdays = 'weekday'
5   #Concatenate pandas objects along a particular axis with optional set logic along the other axes
6   data = pd.concat([Bike['cnt'], Bike[variable_weekdays]], axis=1)
7   #subplot(): The function returns a figure object and a tuple containing axes objects equal to nrows*ncols
8   f, ax = plt.subplots(figsize=(12, 10))
9   #seaborn. boxplot. Draw a box plot to show distributions with respect to categories
10  fig = sns.boxplot(x=variable_weekdays, y="cnt", data=data)
11  fig.axis(ymin=0, ymax=10000)
12  #Boxplot shows most median value between 4000-5000
13
```

```
1   #Bivariate relationship  holiday is catogaries  data or target variable cnt
2   #box plot to show distributions with respect to categories
3   #box plot 'holiday' with 'CNT'
4   variable_holiday = 'holiday'
5   #Concatenate pandas objects along a particular axis with optional set logic along the other axes
6   data = pd.concat([Bike['cnt'], Bike[variable_holiday]], axis=1)
7   #subplot(): The function returns a figure object and a tuple containing axes objects equal to nrows*ncols
8   f, ax = plt.subplots(figsize=(12, 10))
9   #seaborn. boxplot. Draw a box plot to show distributions with respect to categories
10  fig = sns.boxplot(x=variable_holiday, y="cnt", data=data)
11  fig.axis(ymin=0, ymax=10000)
12  # If day is neither weekend nor holiday is 1
13  #otherwise 0
14  #below Boxplot is showing that median  high on holidays when compare to weekdays
```

## Missing value Analysis

```
1   #missing values, occur when no data value is stored for the variable in an observation. Missing data are a common occurrence
2   #and can have a significant effect on the conclusions that can be drawn from the data.
3   #The isnull() function is used to detect missing values for an array-like object
4   #The sort() method sorts the elements of a given list in a specific order - Ascending or Descending.
5   #here we short descending order
6   total = Bike.isnull().sum().sort_values(ascending=False)
7   percent = (Bike.isnull().sum()/Bike.isnull().count()).sort_values(ascending=False)
8   #Total,percent column name
9   missing_data = pd.concat([total, percent], axis=1, keys=['Total', 'Percent'])
10  missing_data.head(100)
11  #In Bike dataset is not available missing value
```

## Outlier Analysis

```
1  # An outlier is an element of a data set that distinctly stands out from the rest of the data.
2  #In other words, outliers are those data points that lie outside
3  #below plot is show variable is available outlier
4  sns.boxplot(Bike['casual'],color="yellow")
```

<matplotlib.axes._subplots.AxesSubplot at 0x2b6f5730e88>

```
1  #in registered numeric variable is not available outlier value
2  sns.boxplot(Bike['registered'],color="yellow")
```

```
1  cnames=['casual']
```

```
1
2  cnames
```

['casual']

```
1   #Detact and delete outlier using of for loop and Iqr
2   for i in cnames:
3       #print variable name
4       print(i)
5       # percentile() function used to compute the nth precentile of the given data (array elements)
6
7       q75,q25=np.percentile(Bike.loc[:,i],[75,25])
8       #IQR can be used to identify outliers in a data set.
9       # Formula of IQR=Q3-Q1
10      iqr=q75-q25
11      min=q25-(iqr*1.5)
12      max=q75+(iqr*1.5)
13      #print minimum value in variable
14      print(min)
15      #print maximum value in variable
16      print(max)
17      #using of drop function we have drope outlier value
18      Bike=Bike.drop(Bike[Bike.loc[:,i]<min].index)
19      Bike=Bike.drop(Bike[Bike.loc[:,i]>max].index)
20      #Below show minimum value in casual variable is -855.25
21      #Below show maximum value in casual variable is 2266.75
22
```

```
3  #After remove outlier in casual variable
4  #plot show their is not available now outliers
5  sns.boxplot(Bike['casual'],color="yellow")
```

## Feature selection

```
1  #Feature selection or variable subset selection,is the process of selecting a subset of relevant
2  #features for use in model construction.
3  #correlation plot
4  Bike_corr=Bike
```

```
1  Bike_corr.shape#total observation is 687 and 16 column
```

(687, 16)

```
1   #set the height and width of the plot
2   #subplots is function set height and width
3   #figsize height =15 or width =10
4   f,ax=plt.subplots(figsize=(15,10))
5   #Generate correlation matrix using of corr funtion
6   corr=Bike_corr.corr()
7   #plot using seaborn library
8   #heatmap is a function available in seaborn library help to plot correlation visualization
9   #diverging_palette function to set color of correlation plot
10  sns.heatmap(corr, mask=np.zeros_like(corr, dtype=np.bool), cmap=sns.diverging_palette(220, 10, as_cmap=True),
11            square=True, ax=ax)
```

```
1   #correlation range is -1 to 1
2   #-1 means highly negitive correlated data plot show blue color
3   #1 meand highly positive correlated data plot show red color
4   #our dataset "atemp" or "temp" indepandent variable is highly positive correlated variable
5   #Independent variable "hum" or target variable "cnt" is no relationship
6   #so we drop two variable atemp or hum variable in our Bike dataset
7   #droping highly correlated variable
8   Bike_day = Bike.drop(['atemp','hum'], axis=1)
9
10
11
```

```
1   #Twp correlated varaible is drop now 14 column is available in Bike dataset
2   Bike_day.shape
```

(687, 14)

```
1   #Top observation after drop atemp or hum variable
2   Bike_day.head()
```

## Feature Scaling

```
1   #Normalization
2   #Normalization is used to scale the data of an attribute so that it falls in a smaller range, such as -1.0 to 1.0 or 0.0 to
3   #In our dataset most numeric variable is normalize form so can't require to check outlier analysis thus variable
4   #Normaliz varaible name is temp,atemp,hum ,windspeed thus are
5   Bike_df =  Bike.copy()
```

```
1   #allready in our dataset four variable is normal form "temp","atemp","hum","windspeed"
2   #two variable casual or registered are not normalized form so we convert their value in normal form
3   Bike.shape
```

(687, 16)

```
1   #Normality check in casual variable
2   plt.hist(Bike["casual"],bins="auto")
```

```
1  vnames=["casual","registered"]
```

```
1  vnames
```

['casual', 'registered']

```
1  #Normalisation
2  for i in vnames:
3      print(i)
4      Bike[i] = (Bike[i] - Bike[i].min())/(Bike[i].max() - Bike[i].min())
```

casual
registered

```
1  #We convert casual variable in normal form
2  Bike['casual'].describe()
```

```
1  #registered variable
2  Bike['registered'].describe()
```

```
1  Bike.head()
```

# Model Selection

```
]:  1  #Splitting our dataset into training set and test set
    2  #Ration of dividing data set is 80 or 20
    3  #split dataset into ratio 80 % train dataset or 20 % test dataset
    4  train, test = train_test_split(Bike_day, test_size=0.2)
```

```
]:  1  #Shape of train dataset
    2  #train dataset available 549 observation or 14 variable
    3  train.shape
```

]:  (549, 14)

```
]:  1  #top random observation in train dataset
    2
    3  train.head()
```

```
#Shape of test dataset
#test dataset available 138 observation or 14 variable

test.shape
```

:, 14)

```
#top 5  observation in test dataset
test.head()
```

## Decision Tree Regressor Model

```
1  #Importing Decison Tree Regressor Module from sklearn
2  from sklearn.tree import DecisionTreeRegressor
3  #We take here feature or data set
4  #train data set to store   object train_features_one independant or target variable
5  #We have to store variable in single object train_features_one independant variable
6  train_features = train[['season','yr','mnth','holiday','weekday','weathersit','temp','windspeed','casual','registered']].val
7  #Target variable
8  train_target_feature = train['cnt'].values
9
10 #test data set to store   object test_feature independant  or target variable
11 #We have to store variable in single object test_feature independant variable
12 test_feature = test[['season','yr','mnth','holiday','weekday','weathersit','temp','windspeed','casual','registered']].values
13 #Target variable
14 test_target_feature= test['cnt'].values
```

```
1  #Decsion Tree for regression
2  DT_fit = DecisionTreeRegressor()
3  #fit function Build a decision tree from the training set (X, y).
4  DT_fit = DT_fit.fit(train_features, train_target_feature)
```

```
1
2  DT_fit
```

```
1  #Apply model on test data
2  #predict function Predict class or regression target for X.
3  #predict is a function the predict unseen array in data set
4  DT_predictions = DT_fit.predict(test_feature)
```

```
1  #This is prediction value
2  DT_predictions
```

## Evaluation decison Tree Regressor model

```
1  #Define MAPE function
2  #Calculate Mean absolute percentage error (MAPE)
3  #we have create MAPE function to calculate
4  def MAPE(y_true, y_pred):
5      mape = np.mean(np.abs((y_true - y_pred) / y_true))*100
6      return mape
7
```

```
1  #Evaluate Decission tree regressor  using of MAPE
2  #Mean absolute percentage error |
3  MAPE(test_target_feature, DT_predictions)
4  #Here MAPE  is a measure of prediction accuracy of a forecasting method in statistics
5
```

17.15871570190131

```
1  #Mean squared error
2  #we have import mean_squared_error to calculate MSE
3  from sklearn.metrics import mean_squared_error
4
5
```

```
1
2  #Mean squared error
3  mean_squared_error(test_target_feature,DT_predictions)
4  #Here MSE that is, the average squared difference between the estimated values and the actual value.
5
```

```
1  #RMSE Root mean squared error
2  #math library available sqrt function
3  #sqrt function calculate square root
4  from math import sqrt
5  result=sqrt(mean_squared_error(test_target_feature,DT_predictions))
6  #here RMSE is a frequently used measure of the differences between values predicted
7  #by a model or an estimator and the values observed.
```

```
1  #Here RMSE
2  result
```

## Random Forest

```
1   #Implemantion of Random forest Regressor Model
2   #Importing Random Forest Regressor Module from sklearn.ensemble pakage
3   from sklearn.ensemble import RandomForestRegressor
4   #train data set to store  object train_features_one independant or target variable
5   #We have to store variable in single object train_features_one independant variable
6   train_features=train[['season','yr','mnth','holiday','weekday','weathersit','temp','windspeed','casual','registered']].value
7   #Target variable
8   train_target_feature = train['cnt'].values
9
10  #test data set to store  object test_feature independant  or target variable
11  #We have to store variable in single object test_feature independant variable
12  test_feature=test[['season','yr','mnth','holiday','weekday','weathersit','temp','windspeed','casual','registered']].values
13  #Target variable
14  test_target_feature= test['cnt'].values
15  #Creating a Random Forest Regressor model and fitting it to the training data
16  #A model comprised of many models is called an Ensemble model.
17  #For this model I've chosen 400 trees (n_estimator=400)
18  RF_model= RandomForestRegressor(n_estimators= 400, random_state=200).fit(train_features,train_target_feature)
```

```
1  #RF_model object
2  RF_model
```

```
1  #Prediction of unseen data using or predict function
2  #prediction array value is
3  RF_predictions
```

# Evaluation Random Forest Model

```
1  #Using of MAPE fucntion we calculate Mean absolute percenatage
2  #Evaluate Random forest using  MAPE
3  MAPE(test_target_feature, RF_predictions)
4  #Mean absolute percentage error
```

231.64270531827168

```
1  #Mean squared error
2  #we have import mean_squared_error to calculate MSE
3  from sklearn.metrics import mean_squared_error
4  mean_squared_error(test_target_feature,RF_predictions)
5  #Here MSE
6
```

8611360.08128433

```
1  #RMSE Root mean squared error
2  #math library available sqrt function
3  #sqrt function calculate square root
4  from math import sqrt
5  result=sqrt(mean_squared_error(test_target_feature,RF_predictions))
```

```
1  #RMSE
2  result
3  #Here Root mean squared error RMSE is given below
```

# Linear Regression Model

```
1   #We take here feature or data set
2   #train data set to store  object train_features_one independant or target variable
3   #We have to store variable in single object train_features_one independant variable
4   train_features = train[['season','yr','mnth','holiday','weekday','weathersit','temp','windspeed','casual','registered']].val
5   #Target variable
6   train_target_feature = train['cnt'].values
7
8   #test data set to store  object test_feature independant  or target variable
9   #We have to store variable in single object test_feature independant variable
10  test_feature = test[['season','yr','mnth','holiday','weekday','weathersit','temp','windspeed','casual','registered']].values
11  #Target variable
12  test_target_feature= test['cnt'].values
```

```
1  #statsmodels is a Python module that provides classes and functions for the estimation of many different statistical models,
2  #as well as for conducting statistical tests, and statistical data exploration
3  import statsmodels.api as sm
4  #implement of Linear regression model
5  #OLS stands for Ordinary Least Squares
6  LR_model= sm.OLS(train_target_feature, train_features).fit()
7
```

```
1  #summary of Linear Regression  model
2  LR_model.summary()
3  #Using of summary function we get summary of Linear Regression model
4  #LR_model
```

```
1  #predicting Linear regression model using predict function
2  LR_predictions = LR_model.predict(test_feature)
```

```
1  LR_predictions
```

## Evaluation Linear Regression Model

```python
#Evaluate Linear regression model  using of MAPE
#Mean absolute percentage error
MAPE(test_target_feature, LR_predictions)
#MAPE
```

1.4168440321532624e-13

```python
#Evaluate MSE
#Mean squared error
mean_squared_error(test_target_feature,LR_predictions)
#Here MSE =2.4407868961336697e-23
```

9.18985272183296e-25

```python
#Evaluate RMSE
#RMSE Root mean squared error
#math library available sqrt function
#sqrt function calculate square root
from math import sqrt
result=sqrt(mean_squared_error(test_target_feature,LR_predictions))
```

```python
result
```

9.586371952846895e-13

```python
"""Conclusion for Bike dataset Linear Regression model is a best model and give
   approximately 99.9 percent accuracy for Bike rental count on daily based """
```

**Appendix C-Extra Figure**

Figure show With outliers in casual variable :



Figure show Without outliers in casual variable:



Figure show Relationship between Weekdays and cnt target variable:

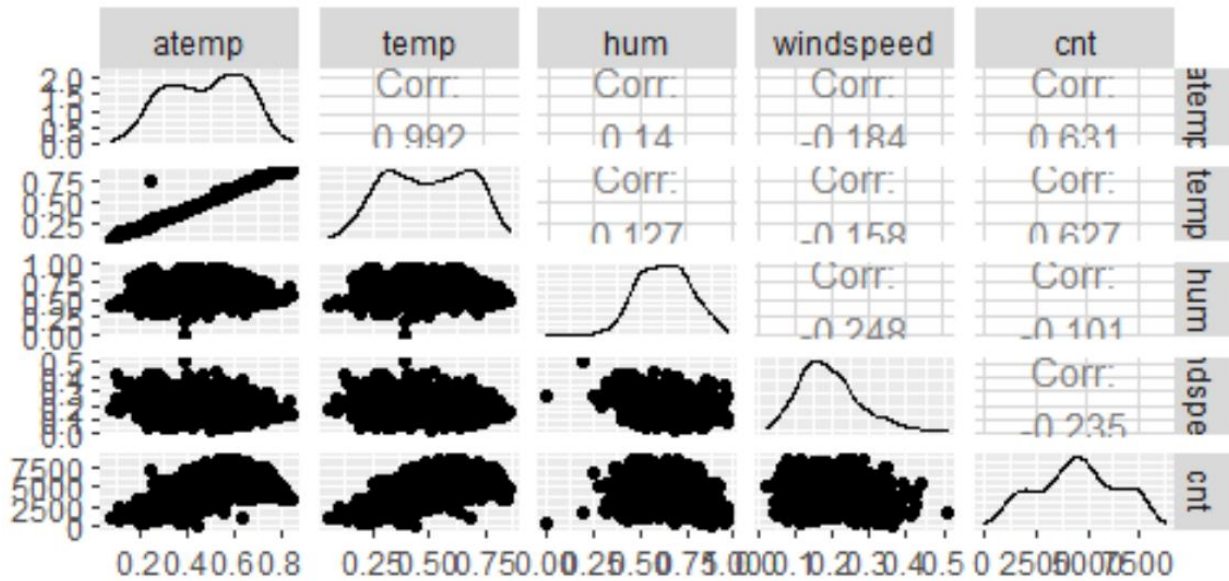Figure show Numeric relationship pairs plot:
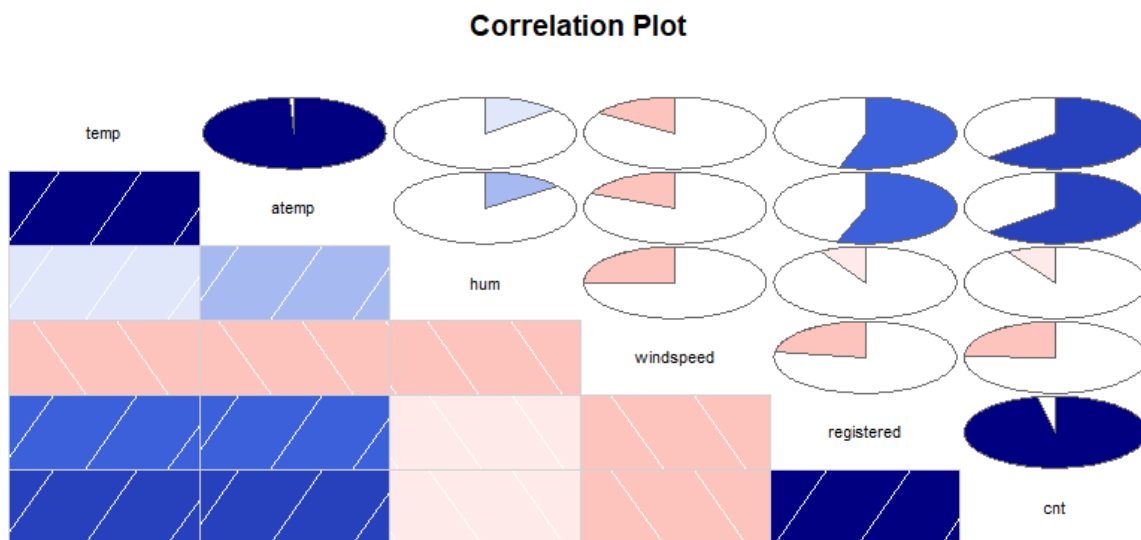


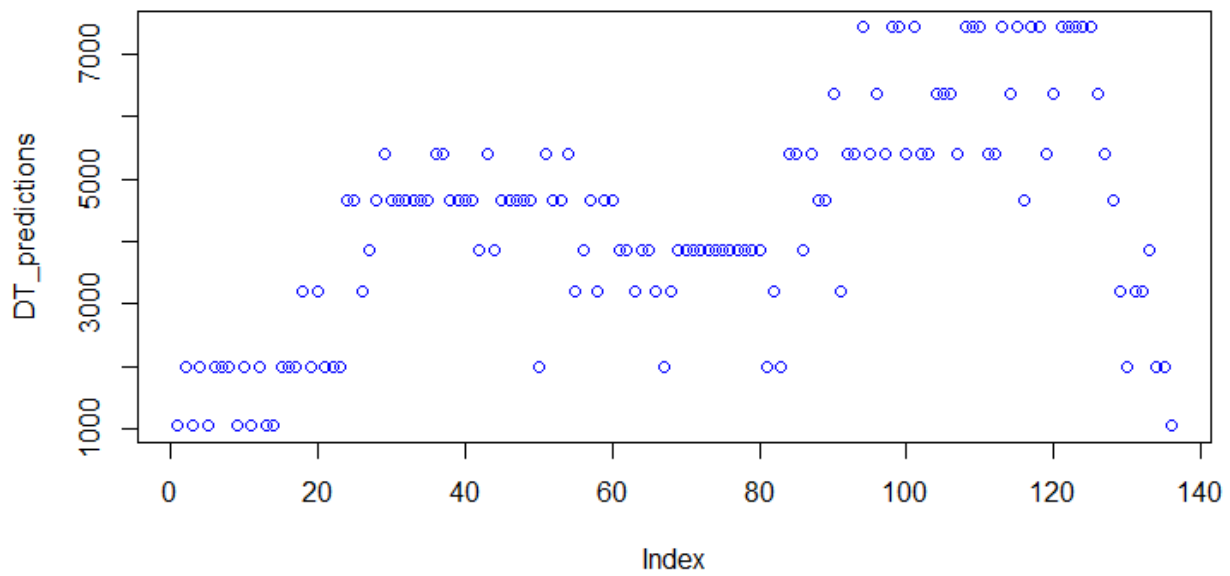Figure show Correlation plot:

Figure shows Decision Tree predictors index:



Figure shows Random Forest error plot:

Figure shows Linear Regression model prediction plot:
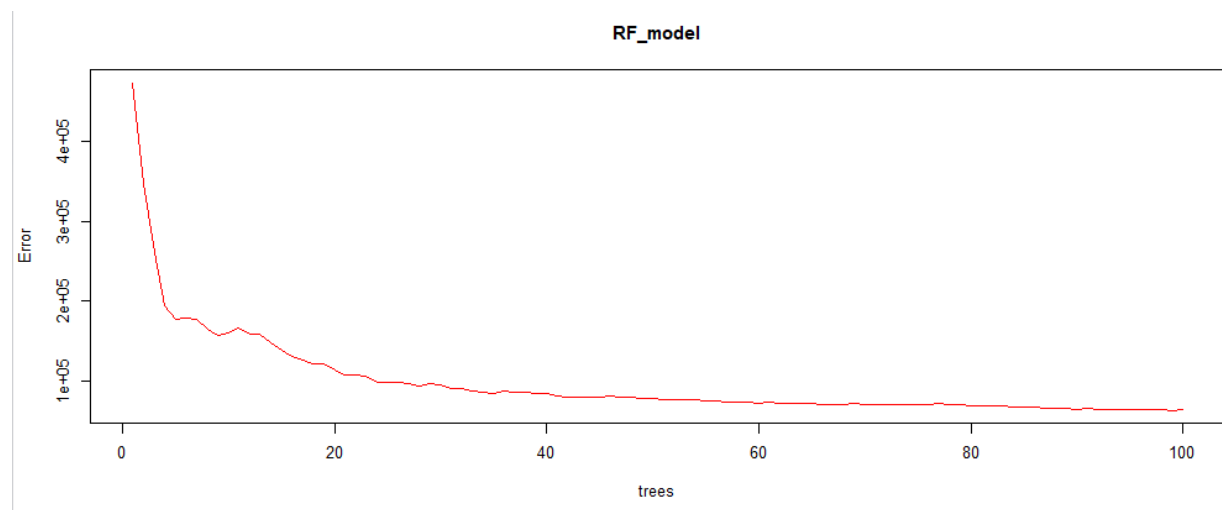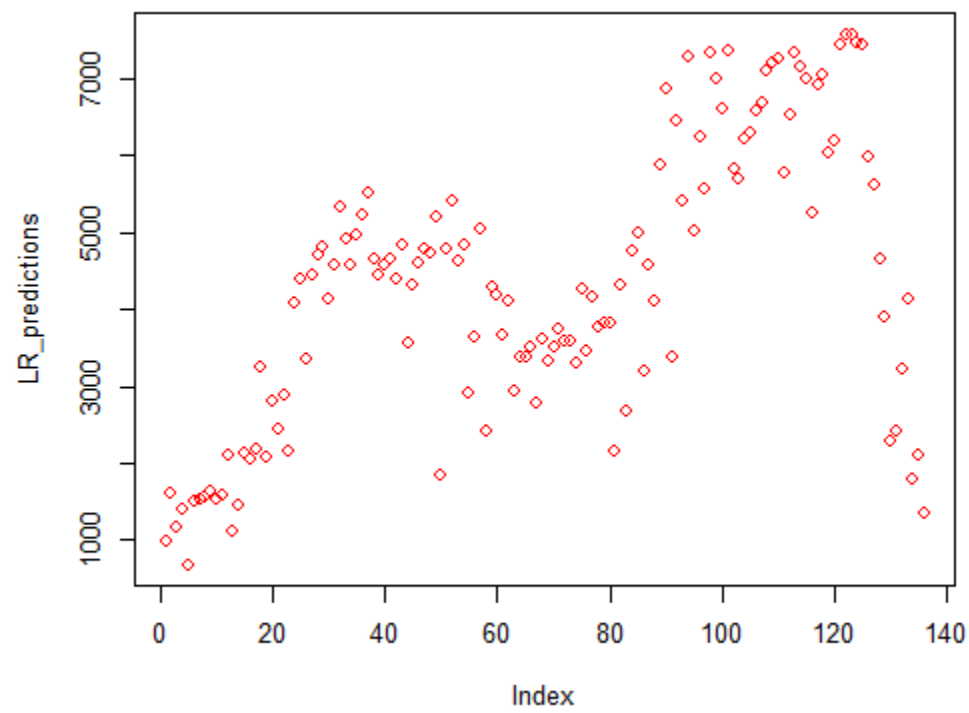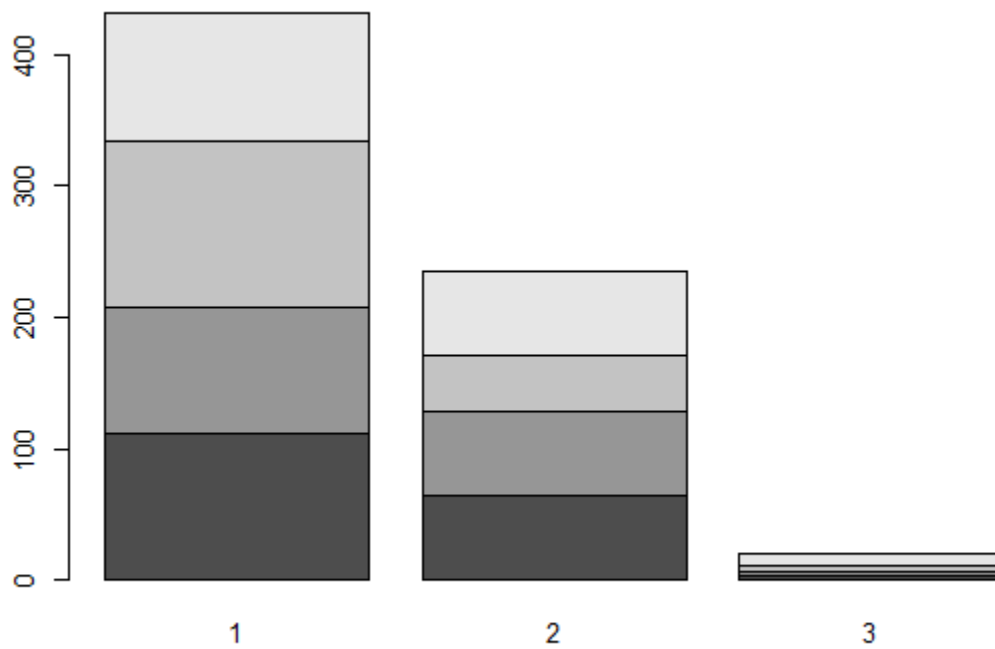


Figure shows Relationship between Holiday and Weathersit:

Figure shows relationship between 'mnth' and 'holiday'