

BIKE PREDICTION

Table of Contents

1.1 Problem Statement

1.2 Dataset

Methodology

2.1 Pre-Processing and EDA

2.1.1 Target Variable – ‘cnt’

2.1.2 Missing value Analysis

2.1.3 Outlier Analysis

2.1.4 Correlation Analysis

2.1.5 Univariate analysis

2.1.6 Bivariate Analysis

2.1.7 Feature Scaling and Normalization

2.2 Modeling

2.2.1 Linear Regression

R Implementation

Python Implementation

2.2.2 Random Forest Regression

R implementation

Python implementation

Result and Performance measure

3.1 Performance Measure

3.1.1 R implementation

3.1.2 Python implementation

3.2 Result

Linear Regression Diagnostics

Chapter 1

1.1 Problem Statement

The objective of this case study is the prediction of bike rental count on daily based on the environmental and seasonal settings. The dataset contains 731 observations, 15 predictors and 1 target variable. The predictors are describing various environment factors and settings like season, humidity etc. We need to build a prediction model to predict estimated count or demand of bikes on a particular day based on the environmental factors.

1.2 Dataset

The data set consist of 731 observation recorded over a period of 2 years, between 2011 and 2012. It has 15 predictors or variables and 1 target variable. All the variables are described in table 1.

Variable names	Description
Instant	Record index
Dteday	Date
Season	Season (1:springer, 2:summer, 3:fall, 4:winter)
Yr	Year (0: 2011, 1:2012)
Mnth	Month (1 to 12)
Hr	Hour (0 to 23)
Holiday	Weather day is holiday or not (extracted from Holiday Schedule)
Weekday	Day of the week
Workingday	If day is neither weekend nor holiday is 1, otherwise is 0.
weathersit	1: Clear, Few clouds, Partly cloudy, Partly cloudy

	2: Mist + Cloudy, Mist + Broken clouds, Mist + Few clouds, Mist 3: Light Snow, Light Rain + Thunderstorm + Scattered clouds, Light Rain + Scattered clouds 4: Heavy Rain + Ice Pallets + Thunderstorm + Mist, Snow + Fog
Temp	Normalized temperature in Celsius. The values are derived via $(t - t_{\min}) / (t_{\max} - t_{\min})$, $t_{\min} = -8$, $t_{\max} = +39$ (only in hourly scale)
Atemp	Normalized feeling temperature in Celsius. The values are derived via $(t - t_{\min}) / (t_{\max} - t_{\min})$, $t_{\min} = -16$, $t_{\max} = +50$ (only in hourly scale)
Hum	Normalized humidity.
Windspeed	Normalized wind speed. The values are divided to 67 (max)
Casual	Count of casual users
Registered	Count of registered users
Cnt	Count of total rental bikes including both casual and registered

Table1. Description of variables

The data set consist of 7 continuous and 8 categorical variables. Sample data is shown below.

instant	dteday	season	yr	mnth	holiday	weekday	workingday	weathersit
1	1/1/2011	1	0	1	0	6	0	2
2	1/2/2011	1	0	1	0	0	0	2
3	1/3/2011	1	0	1	0	1	1	1
4	1/4/2011	1	0	1	0	2	1	1
5	1/5/2011	1	0	1	0	3	1	1
6	1/6/2011	1	0	1	0	4	1	1

temp	atemp	hum	windspeed	casual	registered	cnt
0.344167	0.363625	0.80583 3	0.160446	331	654	985
0.363478	0.353739	0.69608 7	0.248539	131	670	801
0.196364	0.189405	0.43727 3	0.248309	120	1229	1349
0.2	0.212122	0.59043 5	0.160296	108	1454	1562
0.226957	0.22927	0.43695 7	0.1869	82	1518	1600
0.204348	0.233209	0.51826 1	0.0895652	88	1518	1606

Table2. Sample data

Chapter 2

Methodology

The solution of this problem is divided into three parts. First was EDA (Exploratory Data analysis) and pre-processing, followed by modelling and performance tuning and comparison. During first part data pre-processing step like missing value analysis, outlier analysis, univariate and bi-variate analysis etc. were performed. After that data was split into train and test. The target variable is a continuous variable, so it a regression problem. Linear regression and Random forest regression were used for modelling and their performance comparison was performed. Both the algorithms were implemented in R and python.

2.1 Pre-Processing and EDA

Pre-processing and EDA was performed in both R and python. The dataset consists of 731 observations, and 15 predictors. The process of pre-processing and EDA is described below.

2.1.1 Target Variable – ‘cnt’

The target variable in the problem statement is the total count of registered and casual users of bikes on a single day. ‘cnt’ is the combined value of ‘registered’ and ‘casual’ variables. The histogram, distribution and summary statistics of ‘cnt’ are as follow.

Summary Stats	Values
count	731
mean	4504.34
std	1937
min	22
25%	3152
50%	4548
75%	5956
max	8714

Table3. Summary statistics of target variable ‘cnt’

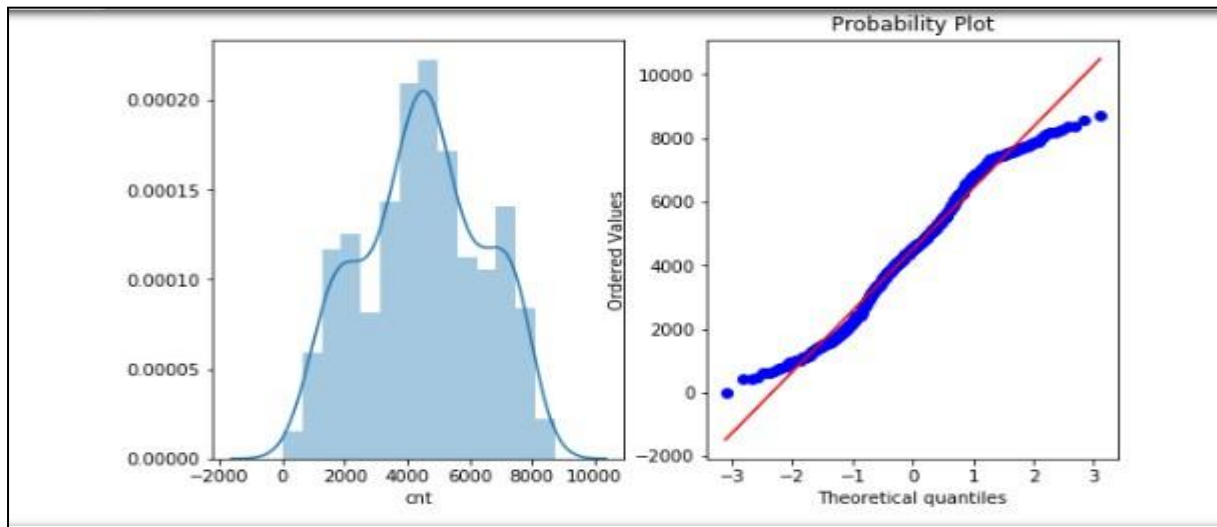


Figure1. Target variable distribution

From the figure it is can be seen that the target variable is close to normal distribution.

2.1.2 Missing value Analysis

Missing value analysis was performed on the dataset. No missing values were found.

2.1.3 Outlier Analysis

After missing value analysis, we check for outliers in target variable and predictors. There were no outliers present in the dataset. Some extreme values were present in the predictors but those seems to logical. So no observations were removed and no imputation was performed on the dataset.

Boxplot method was used to check for outliers. Below are the figures from the python implementation. Box plots from R implementation can be found in appendix.

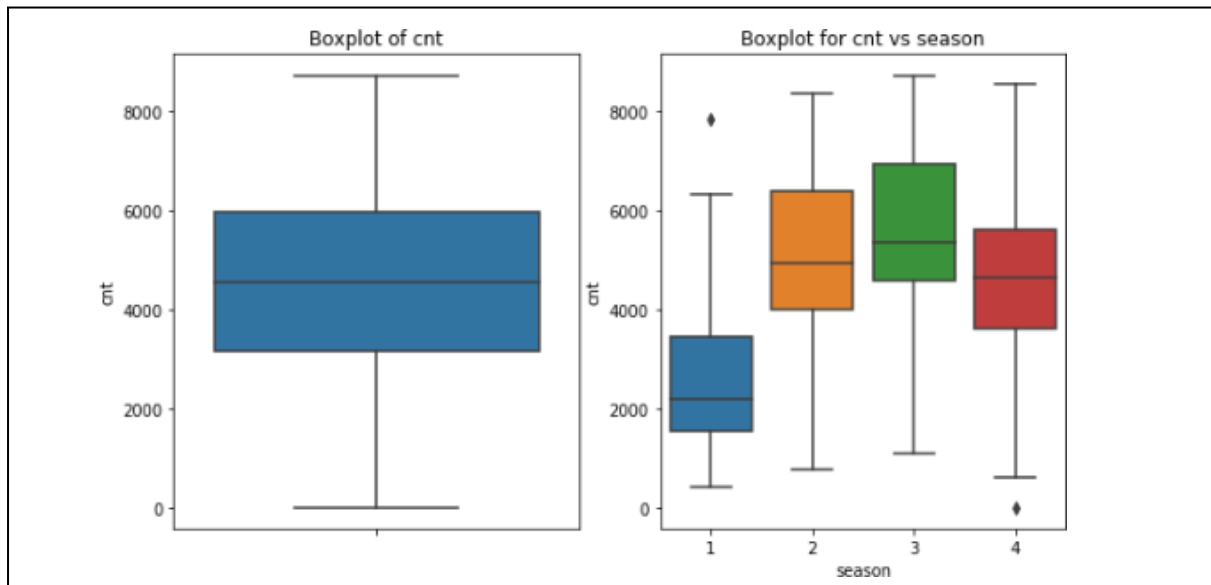


Figure3. Boxplot for cnt and boxplot for cnt vs season

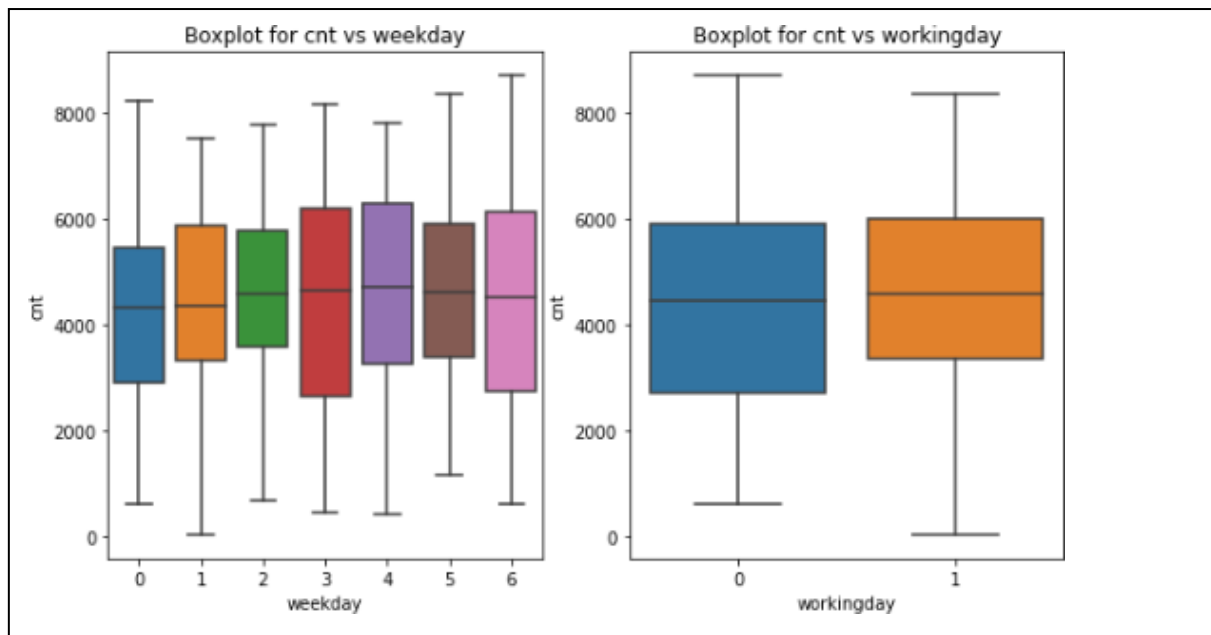


Figure4. Boxplot of cnt vs weekday and boxplot of cnt vs working day

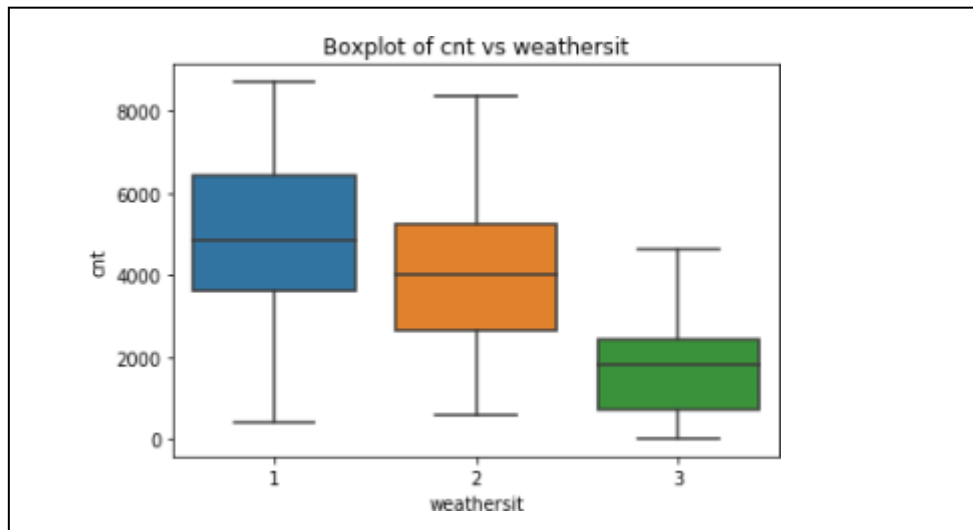


Figure5. Boxplot of cnt vs weathersit

After examining the above boxplots, we can see that there are some extreme values but no outliers. From these boxplots we can also infer that

1. Bike demand count ('cnt') is low in spring (1) season.
2. There is no effect on bike count('cnt') due to a holiday or a working day.
3. Bikes are rented mostly in good weather (1: Clear, Few clouds, Partly cloudy, Partly cloudy) and least in bad (3: Light Snow, Light Rain + Thunderstorm + Scattered clouds, Light Rain + Scattered clouds) weather.

2.1.4 Correlation Analysis

Correlation analysis is used for checking a linear relationship between continuous predictor and target. It is also used to check for multicollinearity among predictors. Multicollinearity exists whenever two or more of the predictors in a regression model are moderately or highly correlated. Multicollinearity is the condition when one predictor can be used to predict other. The basic problem is multicollinearity results in unstable estimation of coefficients which makes it difficult to access the effect of independent variable on dependent variable. Figure6 is showing the correlation matrix for bike rent dataset.

	temp	atemp	hum	windspeed	cnt
temp	1	0.99	0.13	-0.16	0.63
atemp	0.99	1	0.14	-0.18	0.63
hum	0.13	0.14	1	-0.25	-0.1
windspeed	-0.16	-0.18	-0.25	1	-0.23
cnt	0.63	0.63	-0.1	-0.23	1

Figure6. Correlation matrix

‘*registered*’ and ‘*casual*’ were not included in correlation matrix because their sum is equal to the ‘*cnt*’ i.e. Target variable.

From the correlation matrix, it is revealed that

1. Temp (temperature) and atemp (ambient temperature) are highly collinear. One of them should be removed before modeling.
2. ‘cnt’ (demand count) have a strong and positive relationship with temperature and ambient temperature which is logical. People tends to rent bikes more which

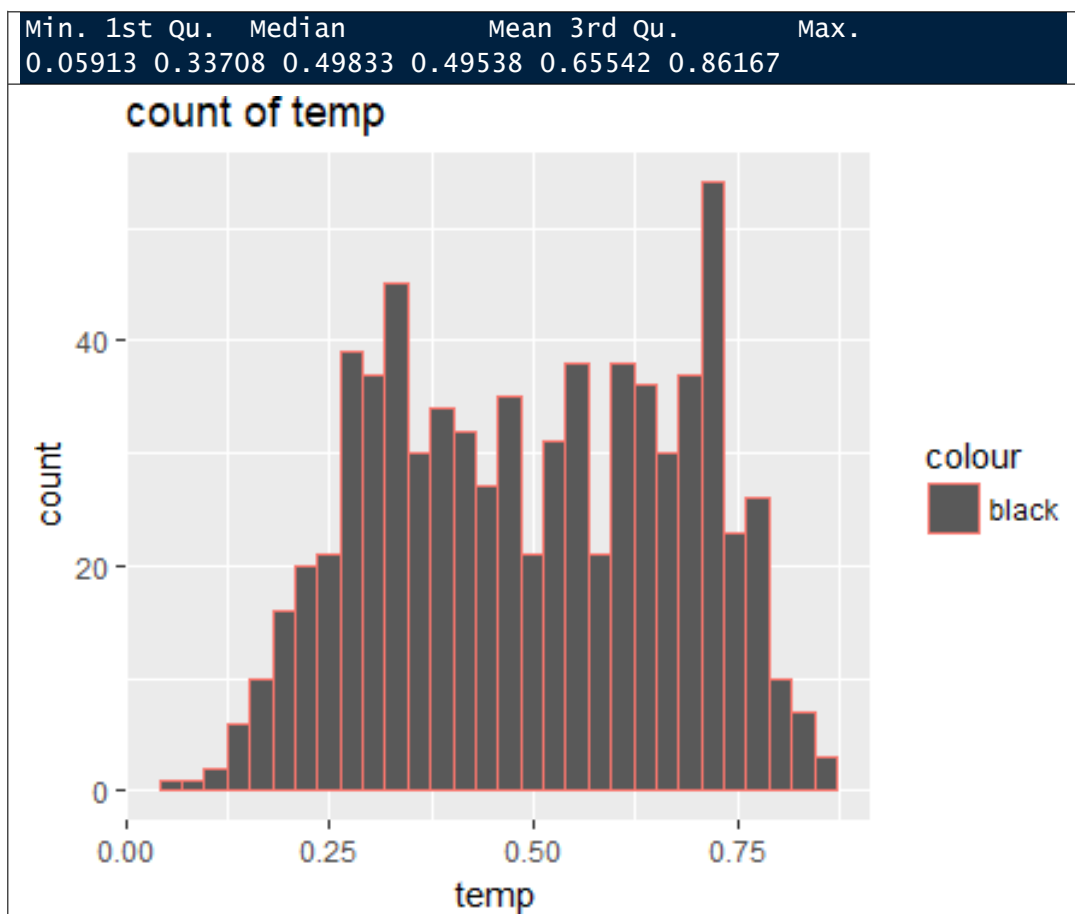
temperature is higher.

3. 'cnt' (demand count) is negative relationship with hum(humidity) and windspeed. People tends to rent bike more when there is less humidity and wind speed.
4. Also the relationship between 'hum', 'windspeed' and 'cnt' is very weak. These are not very strong predictors.

2.1.5 Univariate analysis

In univariate analysis, we look at the distribution and summary statistics of each variable.

- a. temp



- b. Table4. Univariate analysis of temp

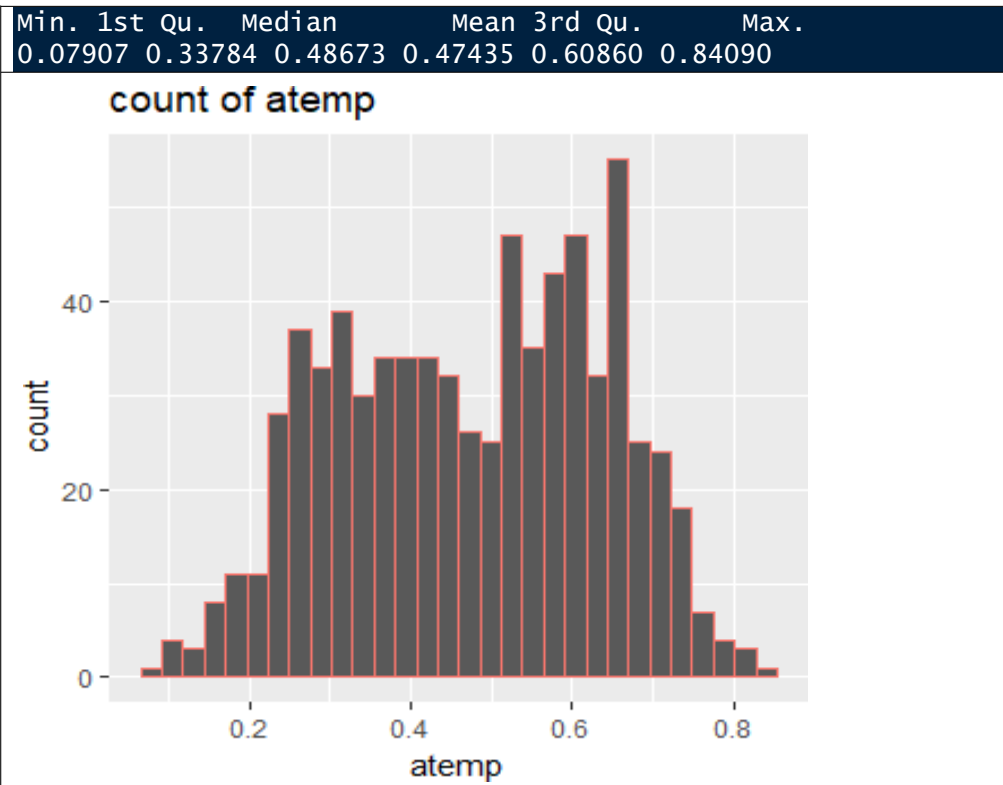


Table5. Univariate analysis of atemp row1. Summary stats row2. Distribution

c. hum

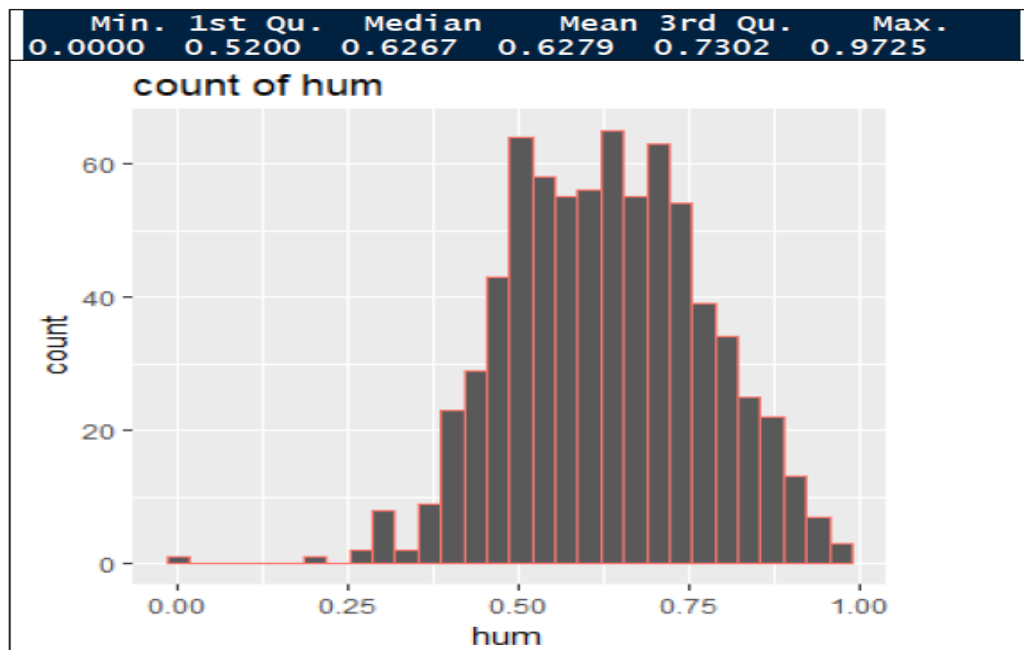


Table6. Univariate analysis of hum(humidity)

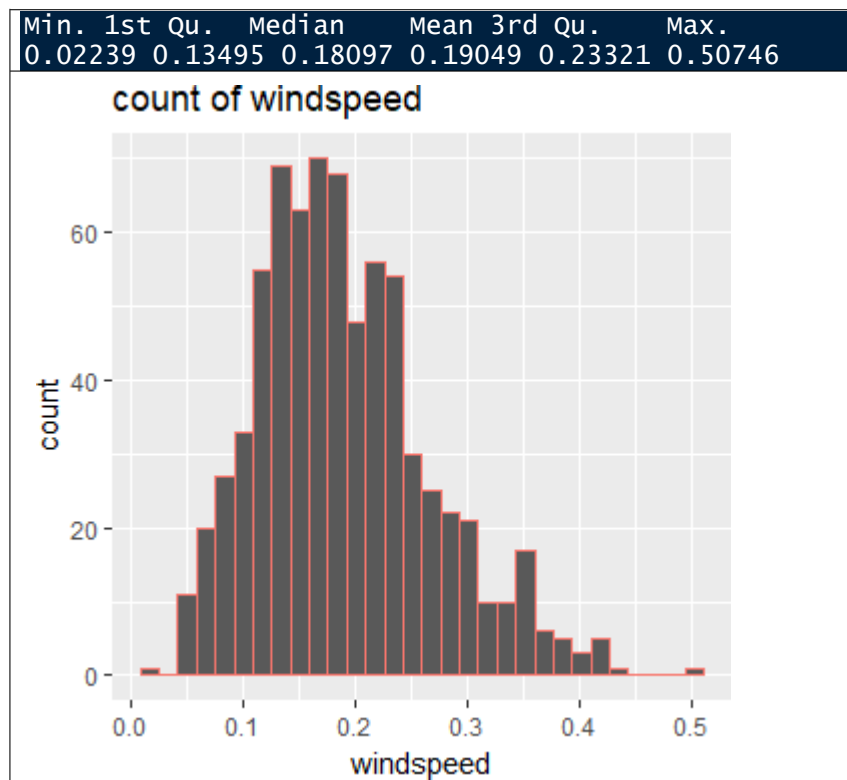


Table7. Univariate analysis of wind summary row1. Summary stats row2. Distribution

d. Seasons

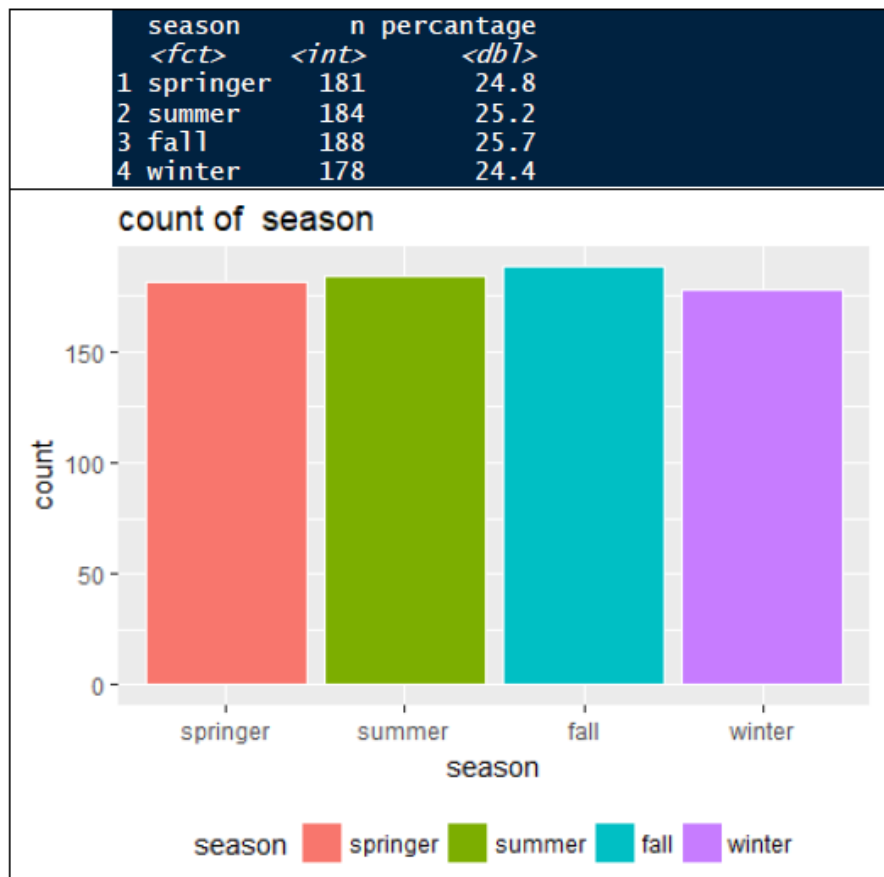


Table8. Univariate analysis of season row1. Summary stats row2. Count

e. Yr(year)

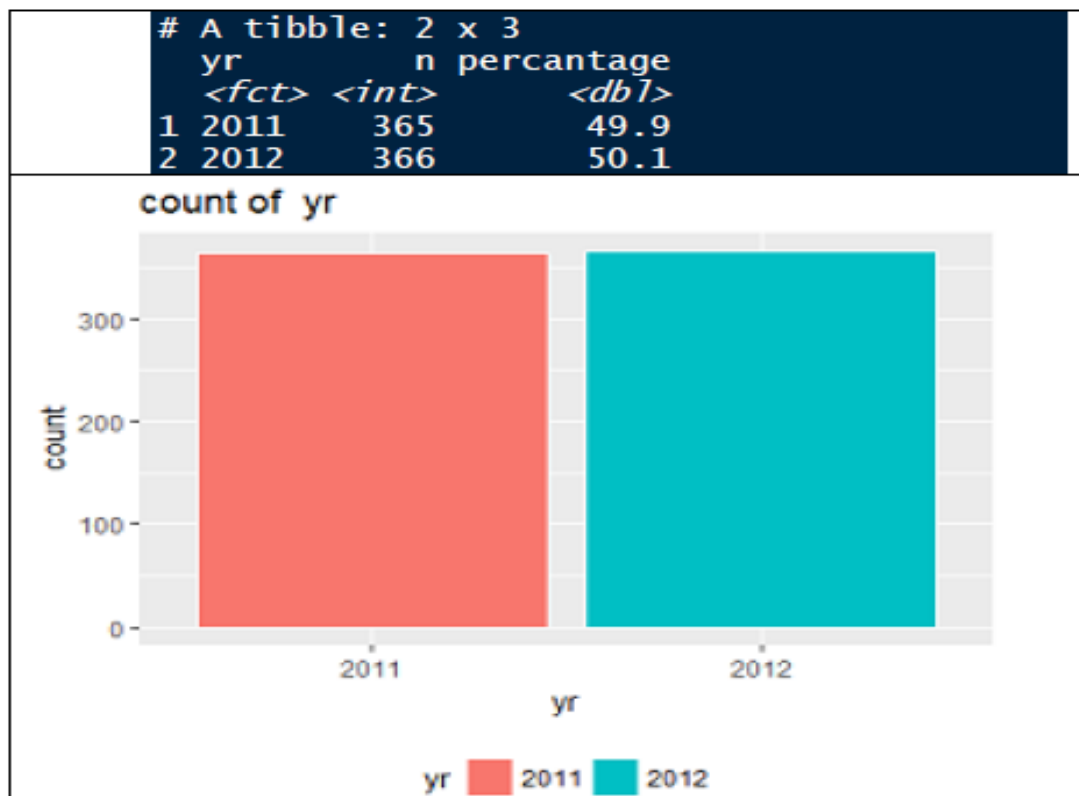


Table9. Univariate analysis of year row1. Summary stats row2. Count

f. Mnth (month)

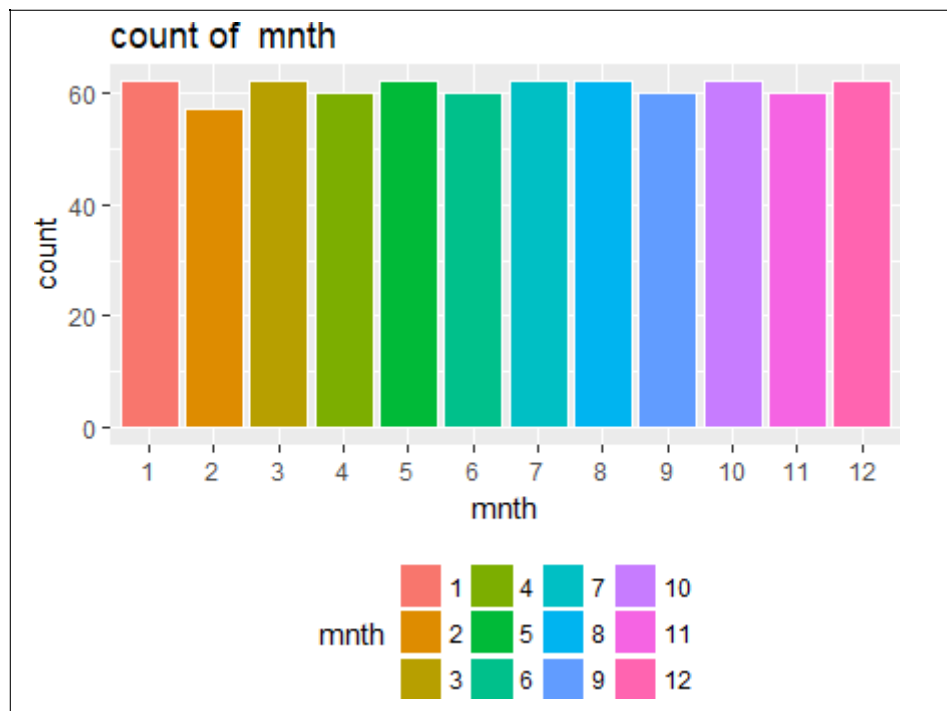


Figure6. count of months

g. Weekday

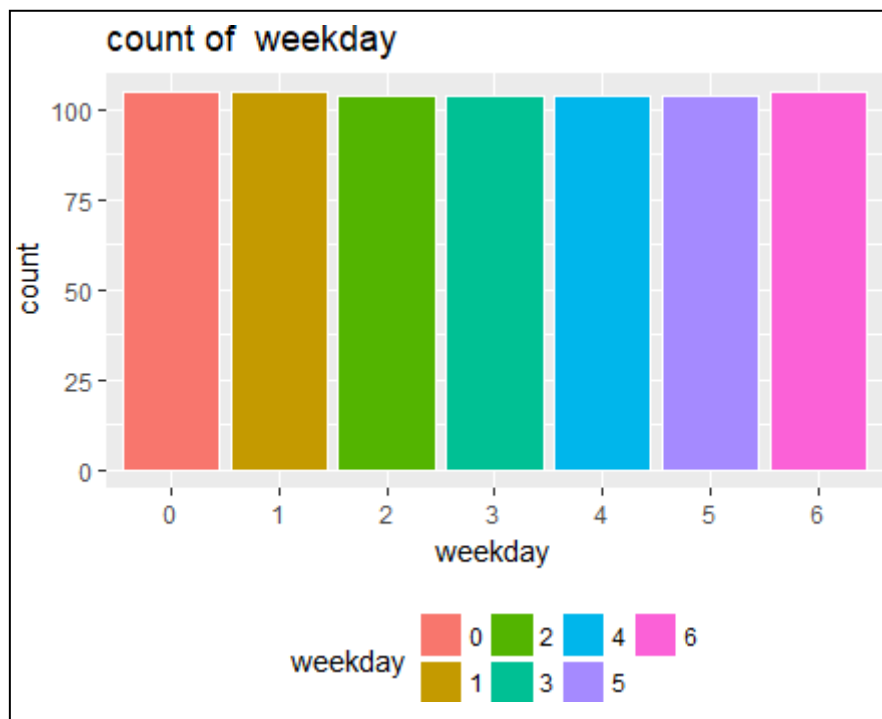


Figure7. count of weekdays

h. Working day

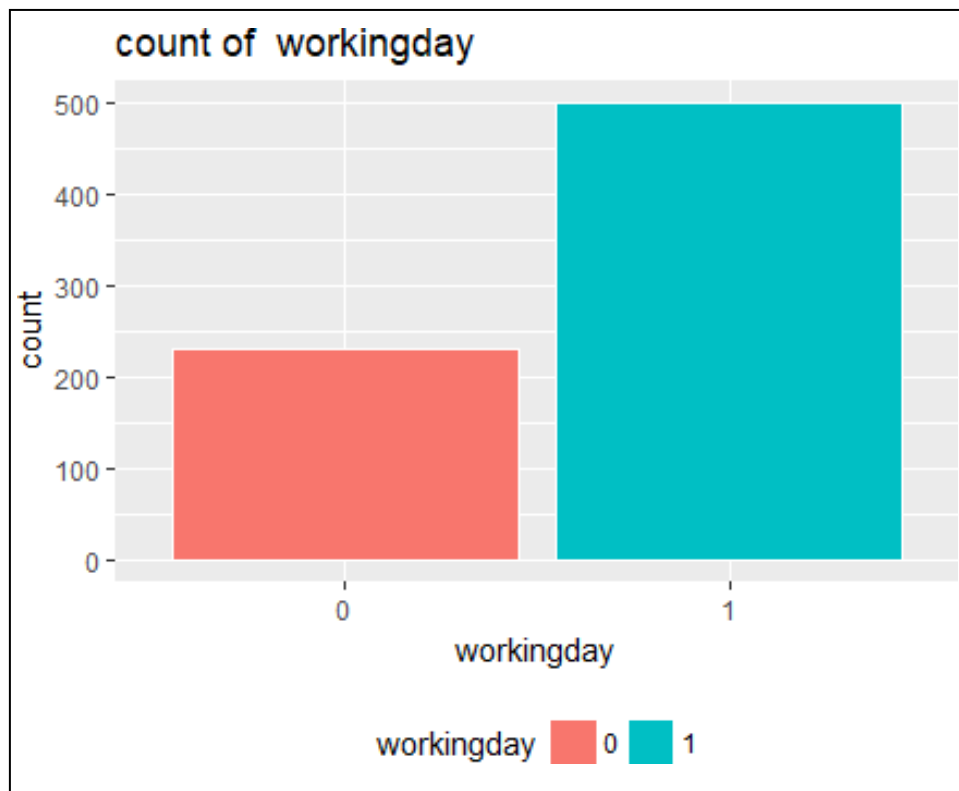


Figure8. count of working days

i. Weathersit (weather situation)

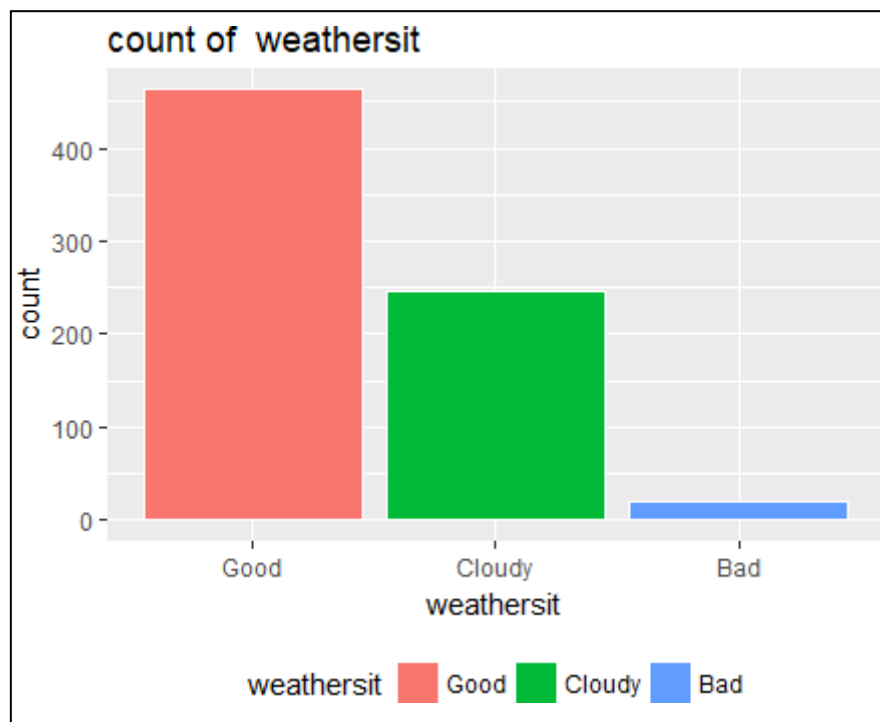


Figure9. count of weather situation

2.1.6 Bivariate Analysis

In bivariate analysis, we will look at the relationship between target variable and predictor. First we look for continuous variables.

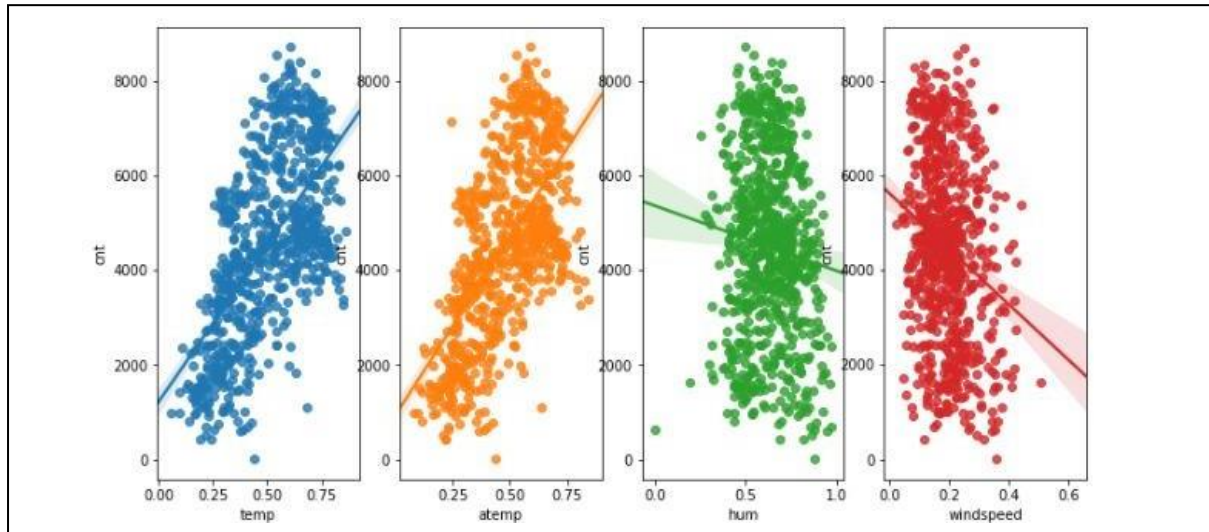


Figure 10. relationship between target variable and continuous predictors From the above scatter plots, we can see that

1. 'cnt' and 'temp' have strong and positive relationship. It means that as the temperature rises, the bike demand also increase.
2. 'atemp' and 'cnt' have strong and positive relationship. It means that as the ambient temperature rise, demand for bikes also increases.
3. 'hum' (humidity) has a negative linear relationship with 'cnt'. As humidity increases, count decreases.
4. 'windspeed' has negative linear relationship with 'cnt'. With an increase in windspeed, bike count decreases.

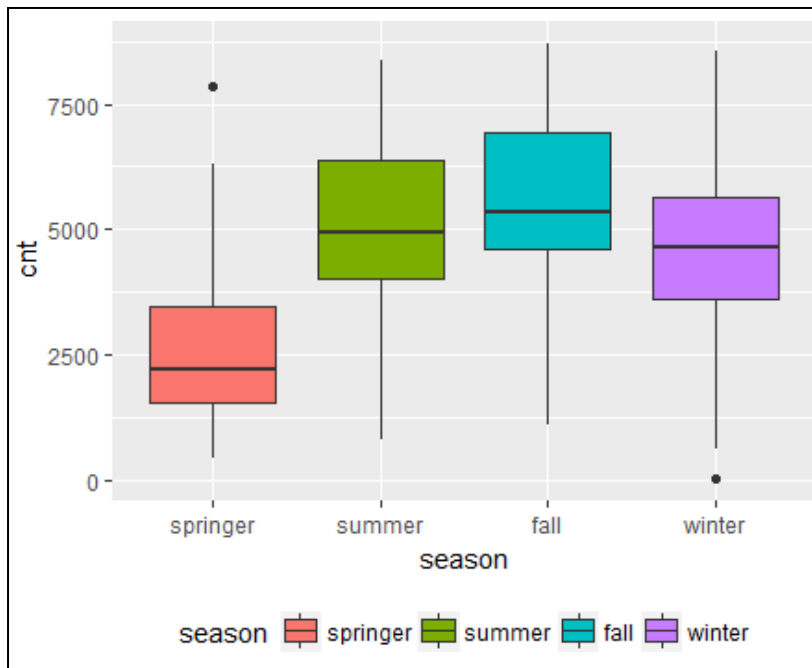


Figure 11 is showing relationship between count (demand) and season.

1. The count is highest for fall season and lowest for springer season.
2. There is no significance difference between count for summer and fall.

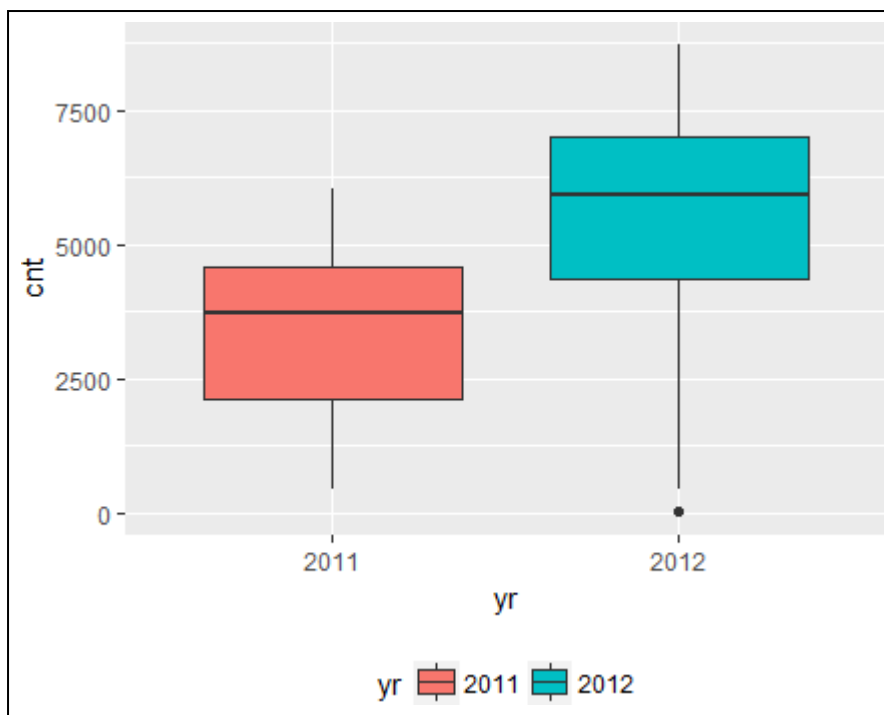
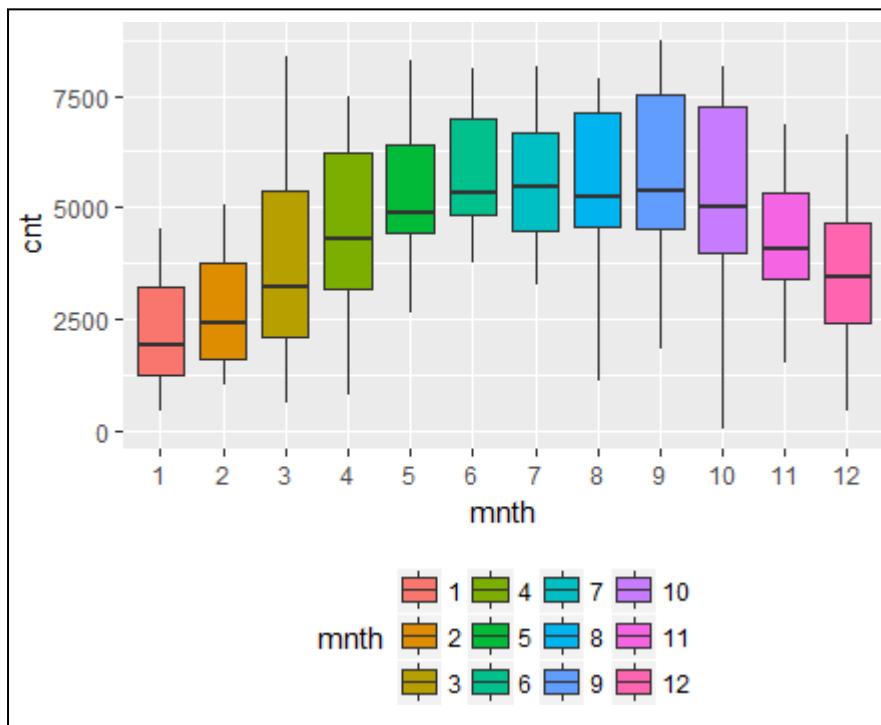


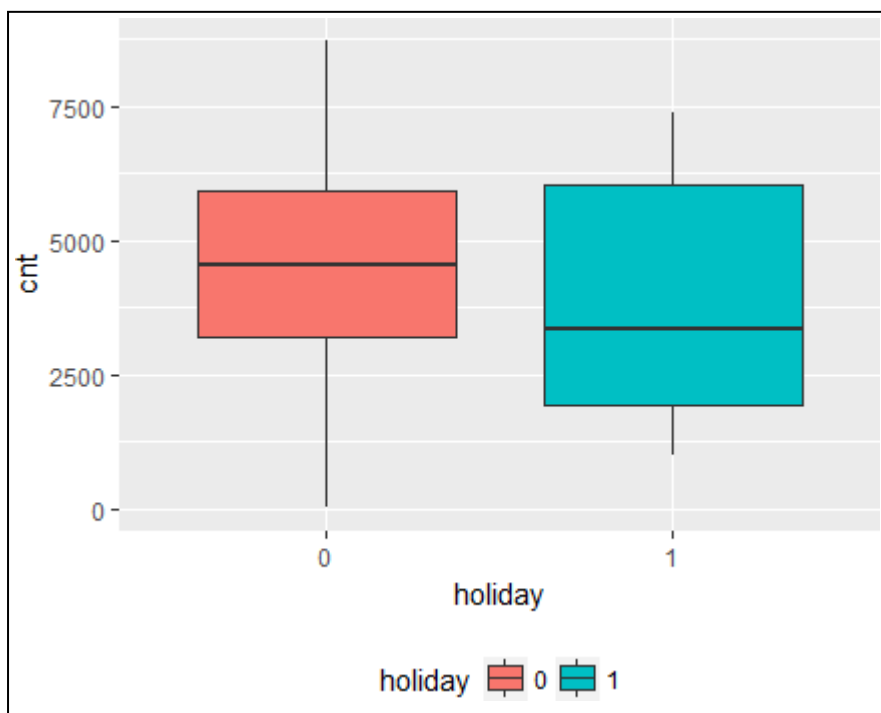
Figure 12 is showing that bike demand was higher in 2012 as compared with 2011.

Figure 12. relationship between year and count



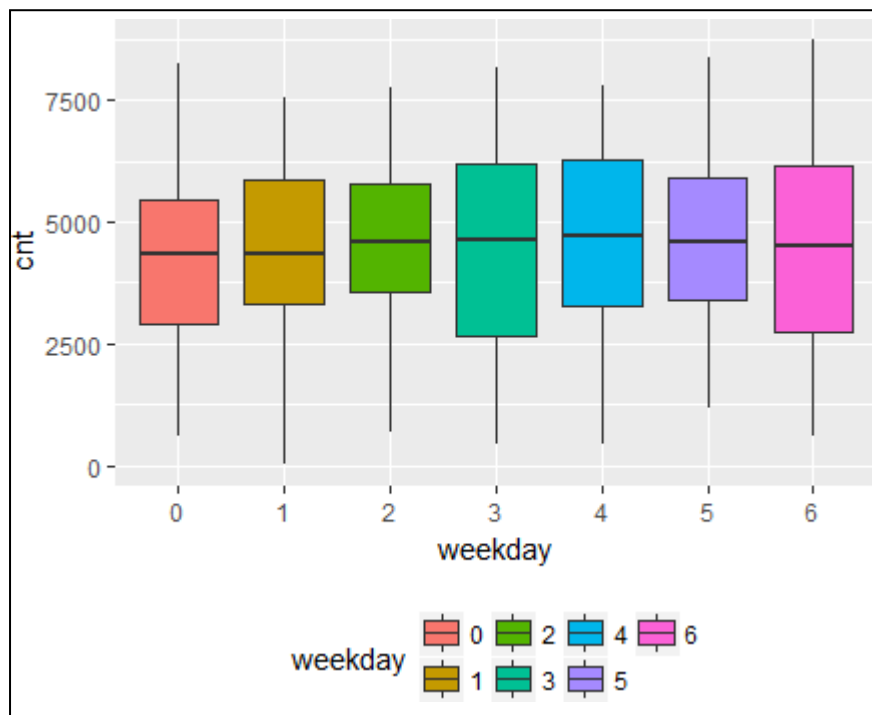
1. From figure 13 it can be inferred that count is high in the month of august, September and October.
2. It is lowest count is for January and February.
3. We can see that as the weather changes from cold to hot, count also

Figure 13. relationship between months and count



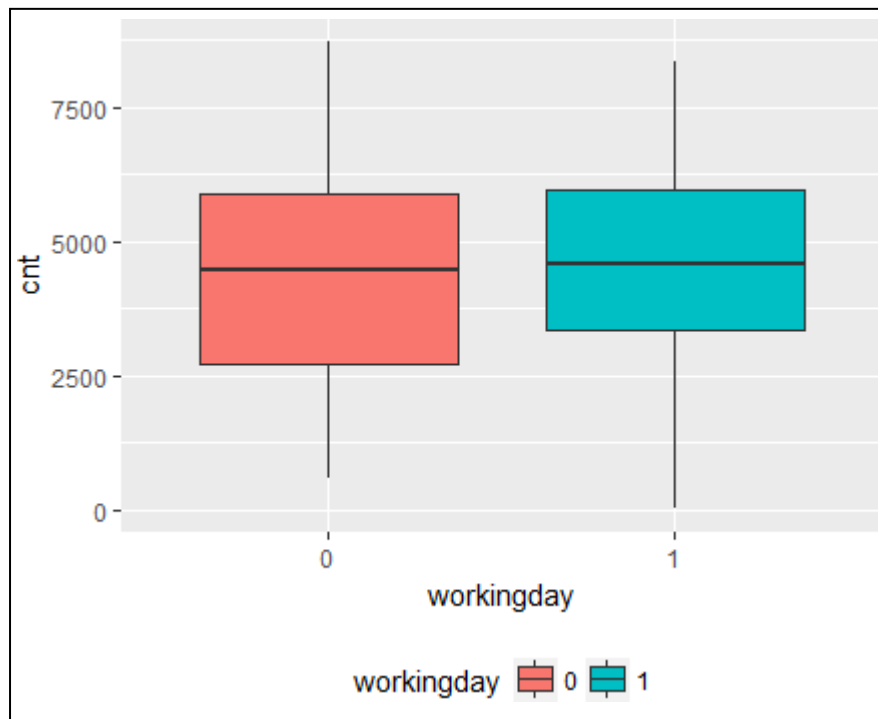
From the boxplot it is visible that count and its median is higher on holidays. People prefer to rent bike on holidays.

Figure 14. relationship between holidays and count



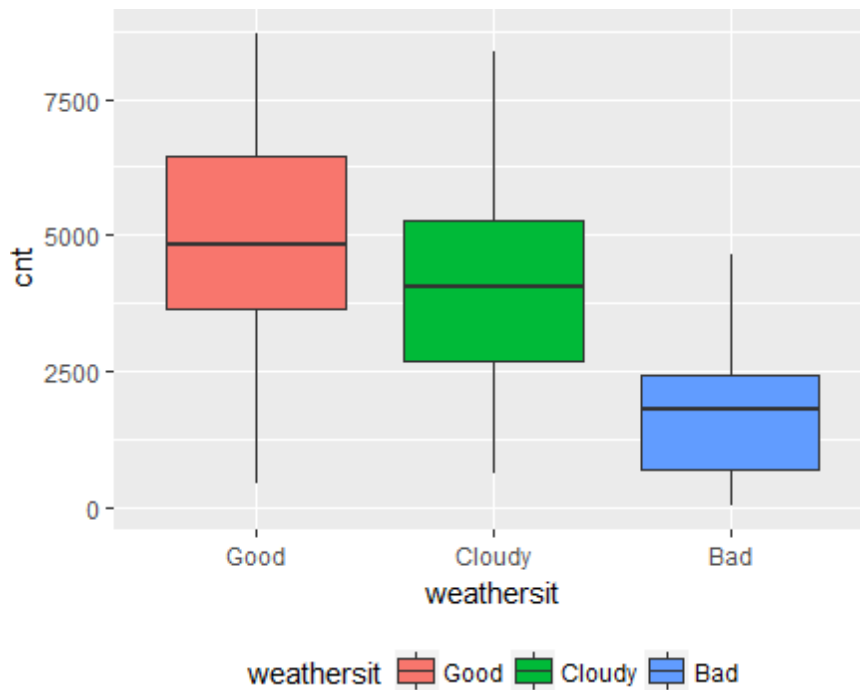
There is not much variation in median of count on weekdays. They are nearly similar on all weekdays.

Figure 15. relationship between weekdays and count



1. There is median for count is same for working and non-working days.
2. The range is longer for non- working days.

Figure 16. relationship between workingday and count



1. The count is maximum when weather situation is good.
2. It is least when weather conditions are bad.

Figure 17. relationship between weathersit and count

2.1.7 Feature Scaling and Normalization

Data normalization is the process of rescaling one or more attributes to the range of $[0, 1]$. This means largest value of each attribute is 1 and smallest is 0. Normalization is a good technique to use when you know that your data distribution is not Gaussian.

Feature scaling was used in the R implementation using MLR package. It was not applied in python for the reason of performance comparison.

2.2 Modeling

In bike renting case study, the target variable is continuous in nature. Our task is predicting the bike demand on a single day. This makes it a regression problem. Two machine learning algorithms were used for learning. Both were implemented in R and python.

1. Multivariate linear regression
2. Random forest regressor – an ensemble tree based regression

After EDA and pre-processing steps, data was divided into training and test dataset with 70 % and 30 % ratio.

After modeling , diagnostic plots were used to check the assumptions of linear regression. For performance tuning of random forest, hyperparameter tuning was used.

2.2.1 Linear Regression

Linear regression is a technique in which we try to model a linear relationship with target and predictors.

First linear regression was used.

- Data was divided into train and test.
- Linear regression was trained on training data.
- Backward and Forward elimination method was used on model with all predictors to select the best model.
- MAP and RMSE was used to check the performance of the model
- Prediction were done on the test data.

R Implementation

First a model with all the predictors was trained in R. I.e. model1. Below is summary of model1. In model1 all the predictors were included. Temp and atemp were multicollinear. They were also included in the season model. The adjusted r square value was .84 which is a good value with F-statistic 106.1

```
1. Call:
2. lm(formula = cnt ~ ., data = training set)
3.
4. Residuals:
5.    Min       1Q   Median       3Q      Max
6. -3518.4  -351.7    63.2   426.9  2439.2
7.
8. Coefficients: (1 not defined because of singularities)
9.              Estimate Std. Error t value Pr(>|t|)
10. (Intercept)    1465.98     296.99   4.936 1.10e-06 ***
11. seasonsummer    1052.92     199.53   5.277 1.99e-07 ***
12. seasonfall     1090.08     243.02   4.486 9.10e-06 ***
13. seasonwinter    1739.77     209.32   8.312 9.69e-16 ***
14. yr2012         2056.81      68.91  29.846 < 2e-16 ***
15. mnth2          207.30      170.46   1.216 0.224539
16. mnth3          505.95      195.71   2.585 0.010026 *
17. mnth4          468.06      284.54   1.645 0.100636
18. mnth5          914.46      310.98   2.941 0.003433 **
19. mnth6          695.86      331.01   2.102 0.036053 *
20. mnth7           74.15      371.79   0.199 0.842003
21. mnth8          537.96      360.50   1.492 0.136280
22. mnth9          954.10      310.59   3.072 0.002247 **
23. mnth10         611.19      285.71   2.139 0.032920 *
24. mnth11        -77.42      268.33  -0.289 0.773079
25. mnth12       -149.70      210.42  -0.711 0.477154
26. holiday1      -477.98      226.35  -2.112 0.035225 *
27. weekday1       84.20      132.11   0.637 0.524209
28. weekday2      216.58      127.58   1.698 0.090235 .
29. weekday3      349.06      126.07   2.769 0.005844 **
30. weekday4      304.14      125.41   2.425 0.015670 *
31. weekday5      374.08      125.54   2.980 0.003030 **
32. weekday6      342.14      124.81   2.741 0.006346 **
33. workingday1    NA          NA          NA      NA
34. weathersitCloudy -409.63      94.17  -4.350 1.66e-05 ***
35. weathersitBad   -2041.38     235.09  -8.684 < 2e-16 ***
36. temp          2548.79     1484.30   1.717 0.086591 .
37. atemp         1571.57     1525.89   1.030 0.303554
38. hum           -1423.48     367.58  -3.873 0.000123 ***
39. windspeed     -2611.79     496.43  -5.261 2.16e-07 ***
40. ---
41. Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
42.
43. Residual standard error: 754.9 on 482 degrees of freedom
44. Multiple R-squared:  0.8561,    Adjusted R-squared:  0.8477
45. F-statistic: 102.4 on 28 and 482 DF,  p-value: < 2.2e-16
```

After model1, step wise model selection was performed. Both forward and backward elimination technique were applied. The second models summary is given below.

```
1. Call:
2. lm(formula = cnt ~ season + yr + mnth + holiday + weekday + weathersit +
3.     temp + hum + windspeed, data = training set)
4.
5. Residuals:
6.    Min       1Q   Median       3Q      Max
7. -3479.9  -351.7    71.3   425.4  2418.5
8.
9. Coefficients:
10.              Estimate Std. Error t value Pr(>|t|)
11. (Intercept)    1514.61     293.24   5.165 3.52e-07 ***
12. seasonsummer    1058.45     199.47   5.306 1.71e-07 ***
13. seasonfall     1092.89     243.02   4.497 8.63e-06 ***
14. seasonwinter    1740.57     209.33   8.315 9.41e-16 ***
15. yr2012         2054.43      68.88  29.826 < 2e-16 ***
16. mnth2          211.07      170.44   1.238 0.216161
17. mnth3          505.08      195.72   2.581 0.010158 *
18. mnth4          471.39      284.54   1.657 0.098240 .
19. mnth5          897.34      310.55   2.889 0.004032 **
20. mnth6          667.54      329.89   2.024 0.043568 *
21. mnth7           53.63      371.28   0.144 0.885217
22. mnth8          488.20      357.27   1.366 0.172427
23. mnth9          928.93      309.64   3.000 0.002839 **
24. mnth10         612.68      285.72   2.144 0.032506 *
25. mnth11        -71.15      268.27  -0.265 0.790960
26. mnth12       -144.34      210.37  -0.686 0.492969
27. holiday1      -493.88      225.83  -2.187 0.029226 *
28. weekday1       84.97      132.12   0.643 0.520427
29. weekday2      212.54      127.53   1.667 0.096254 .
30. weekday3      344.98      126.02   2.738 0.006417 **
31. weekday4      302.66      125.41   2.413 0.016180 *
32. weekday5      365.09      125.24   2.915 0.003721 **
33. weekday6      339.40      124.79   2.720 0.006767 **
34. weathersitCloudy -412.23      94.14  -4.379 1.46e-05 ***
35. weathersitBad   -2050.08     234.47  -8.782 < 2e-16 ***
36. temp          3986.92     503.44   7.919 1.65e-14 ***
37. hum           -1398.37     366.79  -3.812 0.000155 ***
38. windspeed     -2708.02     487.59  -5.554 4.62e-08 ***
39. ---
40. Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
41.
42. Residual standard error: 755 on 483 degrees of freedom
43. Multiple R-squared:  0.8558,    Adjusted R-squared:  0.8477
44. F-statistic: 106.1 on 27 and 483 DF,  p-value: < 2.2e-16
```

Table 11. Summary of AICModel

For second model, adjusted R-square was same, but a slight increase in F-statistics. It shows that this model has better relevant features. There was an issue with this model. If you look at the summary of model, it can be seen that there are negative predictions. We don't want negative values since they don't make any sense. So to counter this, third model was trained with log transformation of target variable. The summary of third model is given below.

```

1. Call:
2. lm(formula = log(cnt) ~ season + yr + mnth + holiday + weathersit +
3.   temp + hum + windspeed, data = training_set)
4.
5. Residuals:
6.   Min       1Q   Median       3Q      Max
7. -4.3758 -0.0998  0.0336  0.1445  0.9327
8.
9. Coefficients:
10.              Estimate Std. Error t value Pr(>|t|)
11. (Intercept)    7.51764    0.11423   65.814 < 2e-16 ***
12. seasonsummer    0.34336    0.08159    4.208 3.06e-05 ***
13. seasonfall      0.41883    0.09910    4.226 2.83e-05 ***
14. seasonwinter    0.64784    0.08562    7.566 1.92e-13 ***
15. yr2012          0.45360    0.02810   16.143 < 2e-16 ***
16. mnth2           0.12485    0.06974    1.790  0.0740 .
17. mnth3           0.13220    0.07991    1.654  0.0987 .
18. mnth4           0.04304    0.11588    0.371  0.7105
19. mnth5           0.09349    0.12656    0.739  0.4604
20. mnth6          -0.09060    0.13407   -0.676  0.4995
21. mnth7          -0.31026    0.15077   -2.058  0.0401 *
22. mnth8          -0.20306    0.14516   -1.399  0.1625
23. mnth9          -0.05415    0.12595   -0.430  0.6675
24. mnth10         -0.12779    0.11680   -1.094  0.2744
25. mnth11         -0.08442    0.10967   -0.770  0.4418
26. mnth12         -0.11101    0.08601   -1.291  0.1974
27. holiday1       -0.20572    0.08900   -2.311  0.0212 *
28. weathersitCloudy -0.04933    0.03825   -1.290  0.1978

29. weathersitBad    -0.90159    0.09560   -9.431 < 2e-16 ***
30. temp            1.70632    0.20373    8.375 5.86e-16 ***
31. hum             -0.59965    0.14895   -4.026 6.58e-05 ***
32. windspeed      -0.96584    0.19933   -4.845 1.70e-06 ***
33. ---
34. Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
35.
36. Residual standard error: 0.3092 on 489 degrees of freedom
37. Multiple R-squared:  0.7364,    Adjusted R-squared:  0.7251
38. F-statistic: 65.04 on 21 and 489 DF,  p-value: < 2.2e-16
39.
40. > summary(predict_test_nonlog)
41.   Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
42.   486   3063   4381   4484   5822  10614

```

Table 12. Summary of stepwiseLogAICModel

In stepwiseLogAICModel model, we lost some adjusted R-square. But after prediction, when the values were converted back using exponential transformation, we get all the positive prediction.

The summary of predictions are as below

```

1. > summary(predict_test_nonlog)
2.   Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
3.   486   3063   4381   4484   5822  10614

```

Chapter 3

Python Implementation

2.2.2 Random Forest Regression

After linear regression, random forest was trained. It was implemented in both R and python. After training with default setting, hyperparameter tuning was used for increase performance.

R implementation

First random forest model was trained 'rf_model_1' with default setting.

```
1. Call:
2. randomForest(formula = cnt ~ ., data = training_set, ntree = 500, mtry = 8, importance = TRUE)
3.      Type of random forest: regression
4.      Number of trees: 500
5. No. of variables tried at each split: 8
6.
7.      Mean of squared residuals: 452950.4
8.      % Var explained: 87.87
```

After hyper tuning , a very litterimprovement was recoreded on variance explained.

```
1. Call:
2. randomForest(formula = cnt ~ . - atemp, data = training_set, ntree = 250, mtry = 6, importance = TRUE)
3.      Type of random forest: regression
4.      Number of trees: 250
5. No. of variables tried at each split: 6
6.
7.      Mean of squared residuals: 451228.2
8.      % Var explained: 87.92
```

In python random forest was trained and hyperparameters optimisation was done using following parameters.

Default setting of random forest are given below.

```
1. RandomForestRegressor(bootstrap=True, criterion='mse', max_depth=None,
2.                        max_features='auto', max_leaf_nodes=None,
3.                        min_impurity_decrease=0.0, min_impurity_split=None,
4.                        min_samples_leaf=1, min_samples_split=2,
5.                        min_weight_fraction_leaf=0.0, n_estimators=10, n_jobs=1,
6.                        oob_score=False, random_state=12345, verbose=0,
7.                        warm_start=False)
```

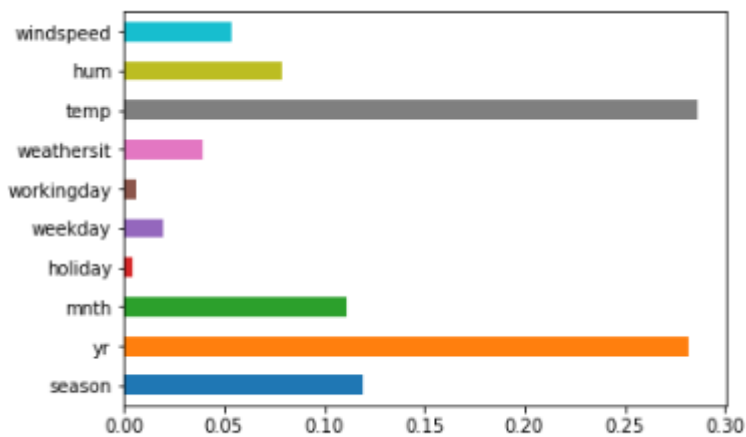
Tuned parameters were selected with hyper tuning random grid search. Best parameters selected were as follows

```
1. Best Parameters using random search:
2. {'n_estimators': 300, 'max_features': 'log2', 'max_depth': 8, 'bootstrap': False}
3. Time taken in random search: 105.77
```

Using above mention parameters, random forest regressor was trained.

```
1. RandomForestRegressor(bootstrap=True, criterion='mse', max_depth=8,
2.                        max_features='log2', max_leaf_nodes=None,
3.                        min_impurity_decrease=0.0, min_impurity_split=None,
4.                        min_samples_leaf=1, min_samples_split=2,
5.                        min_weight_fraction_leaf=0.0, n_estimators=300, n_jobs=1,
6.                        oob_score=False, random_state=12345, verbose=0,
7.                        warm_start=False)
```


Variable importance using random forest in python.



b

Result and Performance measure

RMSE (root mean square error) and MAE (mean absolute error) were used as error metric and measuring model performance.

3.1 Performance Measure

3.1.1 R implementation

For measuring rmse, Metric package was used. For measuring MAE, a function was written. The values for both the metric for linear regression and random forest are as follow.

Error Metric / Algorithm	Linear Regression	Random Forest
RMSE	821.37	745.69
MAE	696.18	498.29

As from the table we can see that random forest performing better than linear regression on both the error metric.

3.1.2 Python implementation

In python, both the error metric was calculated using python functions. No pre-built package or modules were used.

The values for both metric are given below.

Error Metric / Algorithm	Linear Regression	Random Forest
RMSE	1222.15	649.74
MAE	899.5	493.33

As we can see random forest performing better than linear regression.

3.2 Result

From the error metric we can see that random forest is performing better than linear regression in both implementations. The result for random forest is similar in both R and python. But in case of linear regression, R's implementation is performing better than python. The difference here is that data in R was normalized before regression. Selection of model depends on use case. If we want to study the effects of predictors in details, we will go for linear regression and look at the regression equation. If we are care about more precise prediction, we will opt for random forest.