# Energy Forecast Analysis

## Team 4

Ashish Dass, Anamika Jha, Shruti Narain

October 28, 2016

## Abstract

*In this assignment we have focused on developing a model designed for power consumption forecasting using multiple linear regressions. The scripts and the modeling have been done using R. The powers consumed for a single account based on different factors are forecasted. The forward, backward and step wise selection of variables are done to reach to an optimum model for prediction. Use of "Prune Trees" are done to get a visualization of the decision rules of dependent variables.*

# 1. Sample Section

A part of the sample data comes from the data collected in 2014 at the Mildred school in the form of Raw-Data1.csv and Rawdata2.csv. Other half of the data has been extracted online form https://www.wunderground.com/weather. The data from both the sources has been combined using R script and the below modification has been done.

## 1.1  Data from csv  files

Power: The power column contains values in 3 different units (kWh, kVARh and Power factor). The expected result was to get the power in kWh so for every hour of every single day kVARh multiplied with power factor gives kWh power consumption. Hence the value obtained from this multiplication summed with the kWh value gives us the total power consumed for that hour.

- After the power consumption has been obtained for every day, the power value given for every 5 minutes is clubbed together to obtain the value for an hour.

Now that we have a column for power consumption of every day for every hour, we have added separate columns for hour, month, day, year, day of week, weekday and peak hour all of them derived from the date and hour column using the below logic:

- Hour: has been assigned 0-23 values starting from 00:00 to 23:59 based on the values hour column we achieved after clubbing every 5 mins of data.
- Month: January to December has been assigned values 1-12 respectively derived from the date column.
- Day: values assigned from 1-31 based on the date column.
- Year: the year value derived from date column.
- Day of Week: Assigned values 0-6 from Sun-Sat respectively to values taken from date column.
- Weekday: value 1 if the day is a weekday else 0 (with use of chron library).
- Peak hour: for 7 am to 7 pm value 1 has been assigned else 0.

## 1.2    Data from wunderground website

The rest of the columns have been retrieved from wunderground website with the use of devtools library and "Ram-N/weatherData" github install. The data that this website provides us for Boston city is not consistent with the format that we need. Hence data cleaning has been performed to make it consistent with our already created date and time columns.

- Some erroneous value in wind_speedMPH column were replaced with 0.
- Humidity columns had N/A in the records which were replaced by 0.
- erroneous values in Wind_Direction were replaced with NA.
- The clean data was then grouped using Date and hour and mean of Humidity, Visibility, Wind Speed, Seal level pressure, Temperature, Wind Direction degrees were calculated using summarize function.
- Top most values for wind direction and Condition were considered.

The two data sets are then merged based on the hour and date columns.

## 2.Multiple-Linear Regression

Of all the columns we have in our SampleOutput.csv we did variable selection based on the forward, backward and step wise approach.

## Backward Selection:

```
install.packages("leaps")
library(leaps)
#### Backward selection
regfit.bwd=regsubsets(power_~.,data=matWeatherFinal3,nvmax=11, method
="backward")
B=summary(regfit.bwd)
names(B)
B
B$rss
B$adjr2
coef(regfit.bwd,6)
```

```
       Source on Save                                    Run        Source
  1  ### Regression (Subset selection)
  2  ### Needed package and datasets
  3
  4  matweatherFinal3=na.omit(matweatherFinal3) # Get rid of NAs
  5  install.packages("leaps")
  6  library(leaps)
  7
  8  #### Backward selection
  9  regfit.bwd=regsubsets(power_~.,data=matweatherFinal3 ,nvmax=11, method="backward")
 10
 11  B=summary(regfit.bwd)
 12  names(B)
 13  B
 14  B$rss
 15  B$adjr2
 16  coef(regfit.bwd,6)
 17
 18  par(mfrow=c(2,2))
 19  plot(B$rss ,xlab="Number of Variables ",ylab="RSS", type="l")
 20  plot(B$adjr2 ,xlab="Number of Variables ", ylab="Adjusted RSq",type="l")
 21  coef(regfit.bwd ,6)
 22
 18:1    (Top Level)                                                    R Script
```

```
Console E:/Sem4-Fall 2016/ADS/Assignment 2 Materials/
 9  ( 1 )
10  ( 1 )  " "            " "           " "
11  ( 1 )  " "            " "           " "
12  ( 1 )  " "            " "           " "
> B$rss
 [1] 462880793 443825151 423655917 402579078 379305392 356678985 336713251 316293857
 [9] 297164183 278849873 265368671 254489824
> B$adjr2
 [1] 0.03777876 0.07728140 0.11910887 0.16283365 0.21113773 0.25810695 0.29955247
 [8] 0.34195168 0.38167740 0.41971580 0.44770440 0.47028288
> coef(regfit.bwd,6)
(Intercept)     hour__8      hour__9     hour__10     hour__11     hour__12     hour__13
   364.8176     262.6524     278.5281     287.2542     291.9627     282.7587     275.0793
```

- The first mean squared value comes from all the variables selected.
- From step 2, one variable is being removed at a time.
- The variables are dropped one by one till we reach our best result.

| | Sno | Mean R squared | Removed | Equation |
|---|---|---|---|---|
| Step1 | 1 | 0.7211 | | power_ ~ hour__ + month__ + day__ +dayofWeek__ + Weekday__ + Wind_Direction + Conditions + Temperature + Dew_PointF + Humidity + Sea_Level_PressureIn + VisibilityMPH + Wind_SpeedMPH + WindDirDegrees + peakOfhour_ |
| Step 2 | 1 | 0.6308 | hour__ | power_ ~ hour__ + month__ + day__ +dayofWeek__ + Wind_Direction + Conditions + Temperature + Dew_PointF + Humidity + Sea_Level_PressureIn + VisibilityMPH + Wind_SpeedMPH + WindDirDegrees + peakOfhour_ |
| | 2 | 0.6797 | month__ | |
| | 3 | 0.7145 | day__ | |
| | 4 | 0.7069 | dayofWeek__ | |
| | 5 | 0.7211 | Weekday__ | |
| | 6 | 0.7171 | Wind_Direction | |
| | 7 | 0.7189 | Conditions | |
| | 8 | 0.721 | Temperature | |
| | 9 | 0.7205 | Dew_PointF | |
| | 10 | 0.7198 | Humidity | |
| | 11 | 0.7166 | Sea_Level_PressureIn | |
| | 12 | 0.7211 | VisibilityMPH | |
| | 13 | 0.721 | Wind_SpeedMPH | |
| | 14 | 0.7211 | WindDirDegrees | |
| | 15 | 0.7211 | peakOfhour | |

| Step 3 | 1 | 0.6308 | hour__ |
|---|---|---|---|
| | 2 | 0.6797 | month__ |
| | 3 | 0.7145 | day__ |
| | 4 | 0.6046 | dayofWeek__ |
| | 5 | 0.7171 | Wind_Direction |
| | 6 | 0.7189 | Conditions |
| | 7 | 0.721 | Temperature |
| | 8 | 0.7205 | Dew_PointF |
| | 9 | 0.7198 | Humidity |
| | 10 | 0.7211 | Sea_Level_PressureIn |
| | 11 | 0.721 | VisibilityMPH |
| | 12 | 0.721 | Wind_SpeedMPH |
| | 13 | 0.7211 | WindDirDegrees |
| | 14 | 0.7211 | peakOfhour_ |
| | | | |
| Step 4 | 1 | 0.6308 | hour__ |
| | 2 | 0.679 | month__ |
| | 3 | 0.7145 | day__ |
| | 4 | 0.6046 | dayofWeek__ |
| | 5 | 0.7171 | Wind_Direction |
| | 6 | 0.7189 | Conditions |
| | 7 | 0.7209 | Temperature |
| | 8 | 0.7204 | Dew_PointF |

## Summary of the most optimum model obtained using backward selection on train data:

```
1   install.packages("car")
2   library(car)
3   samp <- floor(0.75* nrow(matweatherFinal))
4   samp
5   set.seed(123)
6   training.index <- sample(seq_len(nrow(matweatherFinal)), size= samp)
7   train <- matweatherFinal[training.index, ]
8   test <- matweatherFinal[-training.index, ]
9   h<-matweatherFinal
10
11  na.omit(test)
12  na.omit(train)
13  #check sea level and temperature collinearity
14  lm.fit = lm(power_ ~ hour__ + month__ + day__ +dayofWeek__ +
15              Wind_Direction + Temperature + Dew_PointF +
16              Humidity + Sea_Level_PressureIn ,data=train)
17  summary(lm.fit)
18  RegressionOutput <- summary(lm.fit)$coefficients[,1]
19  write.csv(RegressionOutput, file = "RegressionOutput1.csv")
20
21  vif(lm.fit)
22  install.packages("forecast")
24:26   (Top Level) ÷                                          R Script ÷
```

```
Console E:/Sem4-Fall 2016/ADS/Assignment 2 Materials/
wind_Directionse              31.9707      11.2387    2.045 0.004433
Temperature                   -1.8866       1.0855   -1.738 0.082239 .
Dew_PointF                     4.3813       1.1648    3.761 0.000170 ***
Humidity                      -2.9834       0.5135   -5.810 6.53e-09 ***
Sea_Level_PressureIn         -15.5664       8.3306   -1.869 0.061727 .
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 128.2 on 6426 degrees of freedom
  (52 observations deleted due to missingness)
Multiple R-squared:  0.7177,    Adjusted R-squared:  0.7137
F-statistic: 179.5 on 91 and 6426 DF,  p-value: < 2.2e-16
```

## Regression output:

```
     4  samp
     5  set.seed(123)
     6  training.index <- sample(seq_len(nrow(matWeatherFinal)), size= samp)
     7  train <- matWeatherFinal[training.index, ]
     8  test <- matWeatherFinal[-training.index, ]
     9  h<-matWeatherFinal
    10
    11  na.omit(test)
    12  na.omit(train)
    13  #check sea level and temperature collinearity
    14  lm.fit = lm(power_ ~ hour__ + month__ + day__ +dayofweek__ +
    15              Wind_Direction + Temperature + Dew_PointF +
    16              Humidity + Sea_Level_PressureIn ,data=train)
    17  summary(lm.fit)
    18  RegressionOutput <- summary(lm.fit)$coefficients[,1]
    19  write.csv(RegressionOutput, file = "RegressionOutput1.csv")
    20
    21  vif(lm.fit)
    22  install.packages("forecast")
    23  library(forecast)
    24  pred=predict(lm.fit,test)
    25  PerformanceMetrics <- accuracy(pred,train$power_)
```

18:17   (Top Level) ⇕                                                    R Script

**Console** E:/Sem4-Fall 2016/ADS/Assignment 2 Materials/

```
Multiple R-squared:  0.7177,    Adjusted R-squared:  0.7137
F-statistic: 179.5 on 91 and 6426 DF,  p-value: < 2.2e-16

> RegressionOutput <- summary(lm.fit)$coefficients[,1]
> RegressionOutput
    (Intercept)             hour__1             hour__2
    669.1230718           2.6196710          13.4885581
        hour__3             hour__4             hour__5
     10.0119812          17.9064844          37.0014994
        hour__6             hour__7             hour__8
    236.7689732         356.0020129         378.0166407
        hour__9            hour__10            hour__11
    390.5932896         392.3482593         394.3028698
       hour__12            hour__13            hour__14
```

## Performance matrix:

```
    11  na.omit(test)
    12  na.omit(train)
    13  #check sea level and temperature collinearity
    14  lm.fit = lm(power_ ~ hour__ + month__ + day__ +dayofweek__ +
    15              Wind_Direction + Temperature + Dew_PointF +
    16              Humidity + Sea_Level_PressureIn ,data=train)
    17  summary(lm.fit)
    18  RegressionOutput <- summary(lm.fit)$coefficients[,1]
    19  write.csv(RegressionOutput, file = "RegressionOutput1.csv")
    20
    21  vif(lm.fit)
    22  install.packages("forecast")
    23  library(forecast)
    24  pred=predict(lm.fit,test)
    25  PerformanceMetrics <- accuracy(pred,train$power_)
    26
    27  write.csv(t(PerformanceMetrics), file = "PerformanceMetrics1.csv")
    28
    29  lm.fit1 = lm(power_ ~ ., data=train)
    30  summary(lm.fit1)
    31
    32  pred=predict(lm.fit1.test)
```
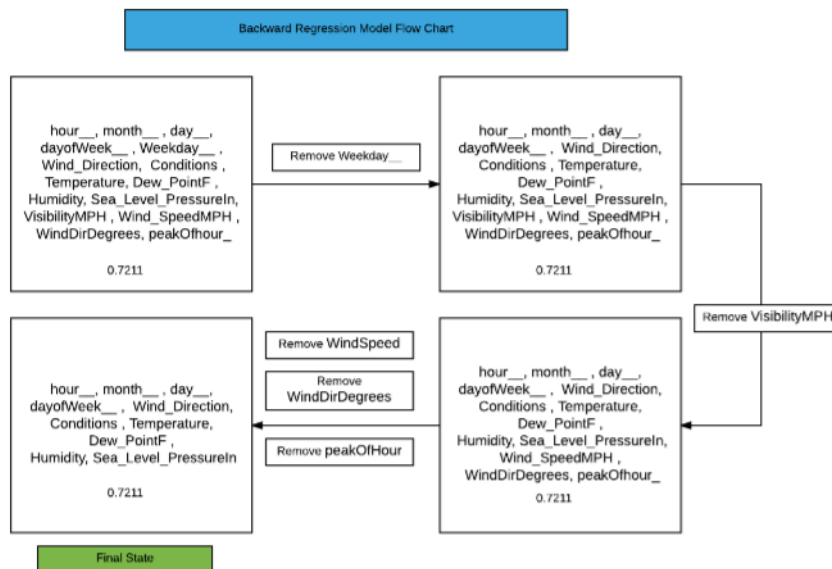
25:9   (Top Level) ⇕                                                    R Script

**Console** E:/Sem4-Fall 2016/ADS/Assignment 2 Materials/

```
   Wind_DirectionEast     Wind_DirectionSSE  Wind_DirectionVariable
           39.8053023            5.1920003             17.2130871
      Wind_DirectionSE           Temperature             Dew_PointF
           31.9707317           -1.8866461              4.3812549
              Humidity   Sea_Level_PressureIn
           -2.9833748          -15.5663799
> t(PerformanceMetrics)
        Test set
ME      -14.29269
RMSE    982.01770
MAE     267.09981
MPE     -31.12927
MAPE     71.48662
>
```

## Flowchart:



## Step wise Selection:

## Summary of the most optimum model obtained using step wise selection on train data:

```r
1   install.packages("car")
2   library(car)
3   samp <- floor(0.75* nrow(matweatherFinal))
4   samp
5   set.seed(123)
6   training.index <- sample(seq_len(nrow(matweatherFinal)), size= samp)
7   train <- matweatherFinal[training.index, ]
8   test <- matweatherFinal[-training.index, ]
9   h<-matweatherFinal
10
11  na.omit(test)
12  na.omit(train)
13  #check sea level and temperature collinearity
14  lm.fit = lm(power_ ~ Temperature +  month__ +  day__ + hour__ +dayofweek__, data=tr
15  summary(lm.fit)
16  RegressionOutput <- summary(lm.fit)$coefficients[,1]
17  write.csv(RegressionOutput, file = "RegressionOutputForward.csv")
18
19  vif(lm.fit)
20  install.packages("forecast")
21  library(forecast)
22  <
```

```
16:1    (Top Level) ÷                                                    R Script
```

```
Console E:/Sem4-Fall 2016/ADS/Assignment 2 Materials/

dayofweek__2 243.1232      0.1210   39.710   < 2e-16
dayofweek__3 232.3008      6.0787   38.215   < 2e-16 ***
dayofweek__4 226.2204      6.0986   37.094   < 2e-16 ***
dayofweek__5 207.5756      6.1180   33.929   < 2e-16 ***
dayofweek__6 104.9973      6.1347   17.115   < 2e-16 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 129.7 on 6446 degrees of freedom
  (52 observations deleted due to missingness)
Multiple R-squared:  0.7102,    Adjusted R-squared:  0.707
F-statistic: 222.5 on 71 and 6446 DF,  p-value: < 2.2e-16

>
```

## Performance Matrix:

```
 8  test <- matWeatherFinal[-training.index, ]
 9  h<-matWeatherFinal
10
11  na.omit(test)
12  na.omit(train)
13  #check sea level and temperature collinearity
14  lm.fit = lm(power_ ~ Temperature +  month__ +  day__ + hour__ +dayofweek__, data=tr
15  summary(lm.fit)
16  RegressionOutput <- summary(lm.fit)$coefficients[,1]
17  write.csv(RegressionOutput, file = "RegressionOutputStepwise.csv")
18
19  vif(lm.fit)
20  install.packages("forecast")
21  library(forecast)
22  pred=predict(lm.fit,test)
23  PerformanceMetrics <- accuracy(pred,train$power_)
24
25  write.csv(t(PerformanceMetrics), file = "PerformanceMetricsForward.csv")
26
27  lm.fit1 = lm(power_ ~ ., data=train)
28  summary(lm.fit1)
29  <
```

21:2   (Top Level) ◊                                                    R Script

Console E:/Sem4-Fall 2016/ADS/Assignment 2 Materials/  ⌂

```
month__      3.732408 11            1.002777
day__        1.324796 30            1.004699
hour__       1.243753 23            1.004753
dayofweek__  1.246360  6            1.018522
> pred=predict(lm.fit,test)
> PerformanceMetrics <- accuracy(pred,train$power_)
> t(PerformanceMetrics)
        Test set
ME      9.321017
RMSE  343.005441
MAE   249.902608
MPE   -27.334111
MAPE   68.733958
```
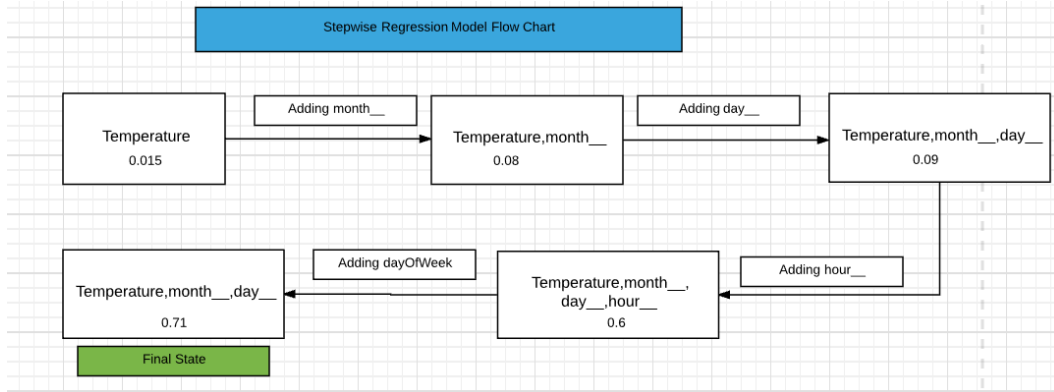
## Regression Output:

```
 4  samp
 5  set.seed(123)
 6  training.index <- sample(seq_len(nrow(matWeatherFinal)), size= samp)
 7  train <- matWeatherFinal[training.index, ]
 8  test <- matWeatherFinal[-training.index, ]
 9  h<-matWeatherFinal
10
11  na.omit(test)
12  na.omit(train)
13  #check sea level and temperature collinearity
14  lm.fit = lm(power_ ~ Temperature +  month__ +  day__ + hour__ +dayofweek__, data=tr
15  summary(lm.fit)
16  RegressionOutput <- summary(lm.fit)$coefficients[,1]
17  write.csv(RegressionOutput, file = "RegressionOutputStepwise.csv")
18
19  vif(lm.fit)
20  install.packages("forecast")
21  library(forecast)
22  pred=predict(lm.fit,test)
23  PerformanceMetrics <- accuracy(pred,train$power_)
24
25  <
```

17:61   (Top Level) ◊                                                   R Script

Console E:/Sem4-Fall 2016/ADS/Assignment 2 Materials/  ⌂

```
Residual standard error: 129.7 on 6446 degrees of freedom
  (52 observations deleted due to missingness)
Multiple R-squared:  0.7102,    Adjusted R-squared:  0.707
F-statistic: 222.5 on 71 and 6446 DF,  p-value: < 2.2e-16

> RegressionOutput <- summary(lm.fit)$coefficients[,1]
> RegressionOutput
 (Intercept)   Temperature      month__2      month__3      month__4      month__5
 -24.4577501     1.7877639   -25.1623248   -22.0503418   -77.0812015   -89.8381607
    month__6      month__7      month__8      month__9     month__10     month__11
 -18.7332802    50.6586782    86.0343823   -46.9085329   -97.1714806   -63.0722209
   month__12        day__2        day__3        day__4        day__5        day__6
 -81.0270548    79.5594709    39.2386569    15.5880244    54.3827346    36.2389357
      day__7        day__8        day__9       day__10       day__11       day__12
```

## Flowchart:



## Forward Selection:

```
install.packages("ISLR")
require(leaps)
install.packages("glmnet")
library(glmnet)
library(ISLR)
regfit.fwd = regsubsets(power_~., data= matWeatherFinal, nvmax=11, method= "forward")
F= summary(regfit.fwd)
names(F)
F
```

- Started with Temperature alone as a dependent variable.
- Repeated the process for all the dependent variables.
- We kept adding one more variable at a time to check any improvement in the summary results.
- These steps were repeated till we reached our best results.

Example 1:

```
fit <- lm(power_ ~ Temperature , data=matweatherFinal)

summary(fit)

Call:
lm(formula = power_ ~ Temperature, data = matweatherFinal)

Residuals:
   Min     1Q Median     3Q    Max
-312.4 -185.9 -113.6  180.1 2423.9

Coefficients:
             Estimate Std. Error t value Pr(>|t|)
(Intercept) 394.66679    4.11112    96.0   <2e-16 ***
Temperature   0.80414    0.06331    12.7   <2e-16 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 236.8 on 8419 degrees of freedom
Multiple R-squared:  0.0188,    Adjusted R-squared:  0.01869
F-statistic: 161.3 on 1 and 8419 DF,  p-value: < 2.2e-16
```

Example 2:

```
fit <- lm(power_ ~ Temperature , data=matweatherFinal)

summary(fit)
```

```
Call:
lm(formula = power_ ~ Dew_PointF, data = matWeatherFinal)

Residuals:
   Min     1Q Median     3Q    Max
-309.7 -186.3 -115.3  180.4 1587.8

Coefficients:
             Estimate Std. Error t value Pr(>|t|)
(Intercept) 414.26357    3.51105 117.989   <2e-16 ***
Dew_PointF    0.55620    0.06254   8.894   <2e-16 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 237.9 on 8419 degrees of freedom
Multiple R-squared:  0.009308,  Adjusted R-squared:  0.009191
F-statistic:  79.1 on 1 and 8419 DF,  p-value: < 2.2e-16
```

**Summary of the most optimum model obtained using forward selection on train data:**

```
 1  install.packages("car")
 2  library(car)
 3  samp <- floor(0.75* nrow(matWeatherFinal))
 4  samp
 5  set.seed(123)
 6  training.index <- sample(seq_len(nrow(matWeatherFinal)), size= samp)
 7  train <- matWeatherFinal[training.index, ]
 8  test <- matWeatherFinal[-training.index, ]
 9  h<-matWeatherFinal
10
11  na.omit(test)
12  na.omit(train)
13  #check sea level and temperature collinearity
14  lm.fit = lm.fit = lm(power_ ~ Temperature +  month__ +  day__ + hour__ +dayofweek_
15  summary(lm.fit)
16  RegressionOutput <- summary(lm.fit)$coefficients[,1]
17  write.csv(RegressionOutput, file = "RegressionOutputForward.csv")
18
19  vif(lm.fit)
20  install.packages("forecast")
21  library(forecast)
22
16:1   (Top Level) ÷                                                    R Script ÷
```

```
Console E:/Sem4-Fall 2016/ADS/Assignment 2 Materials/
dayofweek__3     232.0287    0.0709  36.314   < 2e-16
dayofweek__4     226.6737    6.0962  37.183   < 2e-16 ***
dayofweek__5     207.2720    6.1146  33.898   < 2e-16 ***
dayofweek__6     104.5757    6.1320  17.054   < 2e-16 ***
Sea_Level_PressureIn  23.3449   7.4471   3.135 0.001728 **
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 129.6 on 6445 degrees of freedom
  (52 observations deleted due to missingness)
Multiple R-squared:  0.7106,    Adjusted R-squared:  0.7074
F-statistic: 219.8 on 72 and 6445 DF,  p-value: < 2.2e-16

>
```

## Regression Output:

```
 6  training.index <- sample(seq_len(nrow(matweatherFinal)), size= samp)
 7  train <- matweatherFinal[training.index, ]
 8  test <- matweatherFinal[-training.index, ]
 9  h<-matweatherFinal
10
11  na.omit(test)
12  na.omit(train)
13  #check sea level and temperature collinearity
14  lm.fit = lm.fit = lm(power_ ~ Temperature +  month__ +  day__ + hour__ +dayofweek_
15  summary(lm.fit)
16  RegressionOutput <- summary(lm.fit)$coefficients[,1]
17  write.csv(RegressionOutput, file = "RegressionOutputForward.csv")
18
19  vif(lm.fit)
20  install.packages("forecast")
21  library(forecast)
22  pred=predict(lm.fit,test)
23  PerformanceMetrics <- accuracy(pred,train$power_)
24
25  write.csv(t(PerformanceMetrics), file = "PerformanceMetricsForward.csv")
26
27
```
```
16:17   (Top Level) ÷                                                        R Script
```

```
Console  E:/Sem4-Fall 2016/ADS/Assignment 2 Materials/  ⤶

> RegressionOutput <- summary(lm.fit)$coefficients[,1]
> RegressionOutput
      (Intercept)       Temperature          month__2          month__3
     -733.7741813         2.0649232       -26.2435660       -23.8853171
         month__4          month__5          month__6          month__7
      -84.0064189       -98.7122738       -30.1017513        38.7935028
         month__8          month__9         month__10         month__11
       73.5675632       -60.2127855      -103.5471478       -67.1640057
        month__12            day__2            day__3            day__4
      -86.4487545        81.3099195        41.6924738        17.2645892
           day__5            day__6            day__7            day__8
       56.6081316        38.2805195        47.5385017        50.8267705
           day__9           day__10           day__11           day__12
```

## Performance Matrix:

```
10
11   na.omit(test)
12   na.omit(train)
13   #check sea level and temperature collinearity
14   lm.fit = lm.fit = lm(power_ ~ Temperature +  month__ +  day__ + hour__
15   summary(lm.fit)
16   RegressionOutput <- summary(lm.fit)$coefficients[,1]
17   write.csv(RegressionOutput, file = "RegressionOutputForward.csv")
18
19   vif(lm.fit)
20   install.packages("forecast")
21   library(forecast)
22   pred=predict(lm.fit,test)
23   PerformanceMetrics <- accuracy(pred,train$power_)
24
25   write.csv(t(PerformanceMetrics), file = "PerformanceMetricsForward.csv
26
27   lm.fit1 = lm(power_ ~ ., data=train)
28   summary(lm.fit1)
29
30   pred=predict(lm.fit1,test)
31   <
27:1    (Top Level) ÷
```

```
Console  E:/Sem4-Fall 2016/ADS/Assignment 2 Materials/  ⌂
warning in install.packages :
  package 'forecast' is in use and will not be installed
> library(forecast)
> pred=predict(lm.fit,test)
> PerformanceMetrics <- accuracy(pred,train$power_)
> t(PerformanceMetrics)
        Test set
ME       45.69453
RMSE 1860.27492
MAE    286.47067
MPE    -21.71344
MAPE    74.42619
> write.csv(t(PerformanceMetrics), file = "PerformanceMetricsForward.csv")
```

## 3. Prediction based on forecast data

- The Sample output is generated from Raw datasets merged, after data cleaning and feature selection (after trying Forward, Backward and Stepwise) techniques.
- The feature selection, led us to split the Sample Output into 75% training data set and 25% testing dataset, after trying a possible combinations of 60-40, 80-20, etc.
- Then after fitting the data into the multiple linear regression, we used package "Forecast" to use its predict() function to predict the data on the clean Forecast data missing on the Power/kwh(which was to be predicted).

The summary and the performance matrix of the final model selected using step wise regression

```
21
22 ▾ ################# CHECKING FIT OF 1 - THE BEST MODEL #########
23
24   lm.fit = lm(kwh ~ Temperature + month + day + `Day of week` + hour,data=train) # removed pe
25   summary(lm.fit)
26
27 ▾ ################## Getting all coefficients  #########
28   summary(lm.fit)$coefficients[,]
29   |
30 ▾ ◄                    Ⅲ                                                          ►
```
29:1   ⊞ Getting all coefficients ⇕                                                    R Script ⇕

Console C:/Users/310250154/Desktop/ADS/Assignments/assignment 2/Assignment 2(1)/ ⇱   ─☐

```
4482              Scattered Clouds      240.00000
7251                Mostly Cloudy      310.00000
4532                Mostly Cloudy      110.00000
7420                   Light Rain      100.00000
4811                Mostly Cloudy      100.00000
3467                Mostly Cloudy      220.00000
7774                Partly Cloudy      360.00000
5804                     Overcast      240.00000
3383                     Overcast      135.00000
151                      Overcast      350.00000
 [ reached getOption("max.print") -- omitted 5771 rows ]
> lm.fit = lm(kwh ~ Temperature + month + day + `Day of week` + hour,data=train) # removed peakh
our and weekday #humidity is making a diff of 0.001
> summary(lm.fit)

Call:
lm(formula = kwh ~ Temperature + month + day + `Day of week` +
    hour, data = train)

Residuals:
    Min      1Q  Median      3Q     Max
-572.04  -67.31   -4.41   57.08  737.76
```

```
Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept)   33.80554   13.59021   2.487 0.012890 *
Temperature    0.06962    0.03822   1.822 0.068548 .
month2       -22.52104    7.92725  -2.841 0.004512 **
month3       -11.71157    7.74514  -1.512 0.130553
month4       -48.37865    7.92839  -6.102 1.11e-09 ***
month5       -36.56030    7.73785  -4.725 2.35e-06 ***
```

```
hour20       126.29005   11.21701  11.259  < 2e-16 ***
hour21        34.66418   11.08578   3.127 0.001774 **
hour22        12.87687   11.14207   1.156 0.247847
hour23        11.40202   11.17354   1.020 0.307554
```

```
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 129.2 on 6446 degrees of freedom
  (52 observations deleted due to missingness)
Multiple R-squared:  0.7141,    Adjusted R-squared:  0.711
F-statistic: 226.8 on 71 and 6446 DF,  p-value: < 2.2e-16

>
```

- with the use of vif() [variance inflation factors] function we have determined the collinearity between dependent variables:

- Values more than 5 denotes collinearity between the variables. In our case all values are between 0-1 hence the dependent variables we have selected are no collinear.

```
28   summary(lm.fit)$coefficients[,]
29
30 ▾ ################## Getting 1 coefficient   #########
31   coef(summary(lm.fit))["month__02","Estimate"]
32   library("car")
33   vif(lm.fit)
34
35   pred=predict(lm.fit,test)
36   library("forecast")
37   accuracy(pred,train$kwh)
38
39 ▾ ################# CHECKING FIT OF ALL THE VARIABLES ########
40
41   lm.fit1  =  lm(power_ ~ ., data=train)
42   summary(lm.fit1)
43
44   pred=predict(lm.fit1,test)
45   accuracy(pred,train$kwh)
46
47 ▾ ################## PREDICTING THE FORECAST##################
48
49 ◀                              III
```

35:1    # Getting 1 coefficient ⇕

Console C:/Users/310250154/Desktop/ADS/Assignments/assignment 2/Assignment 2(1)/ ⤳

```
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 129.2 on 6446 degrees of freedom
  (52 observations deleted due to missingness)
Multiple R-squared:  0.7141,     Adjusted R-squared:  0.711
F-statistic: 226.8 on 71 and 6446 DF,   p-value: < 2.2e-16

 vif(lm.fit)
                   GVIF Df GVIF^(1/(2*Df))
Temperature    1.167016  1        1.080285
month          1.210154 11        1.008708
day            1.267092 30        1.003953
Day of week` 1.216142  6        1.016441
hour           1.052748 23        1.001118
```

```
25   summary(lm.fit)
26
27 - ################## Getting all coefficients   #########
28   summary(lm.fit)$coefficients[,]
29
30 - ################## Getting 1 coefficient    #########
31   coef(summary(lm.fit))["month__02","Estimate"]
32   library("car")
33   vif(lm.fit)
34
35   pred=predict(lm.fit,test)
36   library("forecast")
37   accuracy(pred,train$kwh)
38
39 - ################# CHECKING FIT OF ALL THE VARIABLES #########
40
41   lm.fit1 = lm(power_ ~ ., data=train)
42   summary(lm.fit1)
43
44   pred=predict(lm.fit1,test)
45   accuracy(pred,train$kwh)
46
47 - ################## PREDICTING THE FORECAST##########################################################
48
49
50   pred=predict(lm.fit,forecastPrediction)
51
```

39:1    # CHECKING FIT OF ALL THE VARIABLES

Console  C:/Users/310250154/Desktop/ADS/Assignments/assignment 2/Assignment 2(1)/

```
- accuracy(pred,train$kwh)
               ME      RMSE      MAE      MPE     MAPE
Test set 0.5016031 313.4098 246.9187 -29.33255 68.78474
-
```

## Prune Tree:

Prune Tree is used to avoid overfitting the data. The size of the tree can be selected in order to minimize the cross validated error.

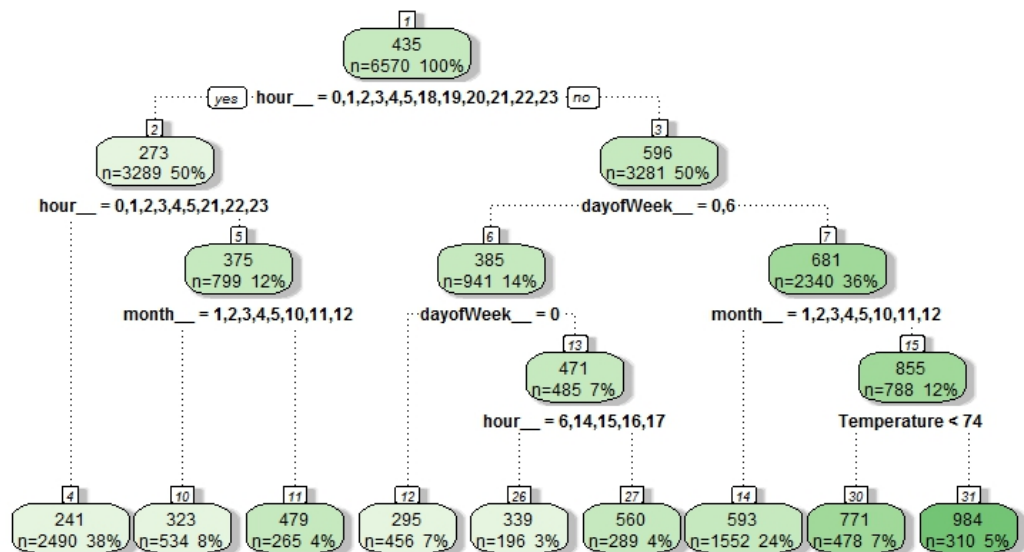- Prune the tree to the desired size using
  prune(fit, cp= )
- printcp( ) is used to examine the cross-validated error results, select the complexity parameter associated with minimum error, and place it into the **prune( )** function

The prune tree results for our model is pasted below:

```
 1  install.packages("rpart.plot")
 2  install.packages("rpart")
 3  library(rpart)
 4  library(rattle)
 5
 6
 7  rt <- rpart(power_ ~ Temperature +  month__ +  day__ + hour__ +dayofWeek__ , data=train)
 8
 9  test.pred.rtree <- predict(rt,test)
10
11
12  printcp(rt)
13  min.xerror <- rt$cptable[which.min(rt$cptable[,"xerror"]),"CP"]
14  min.xerror
15  rt.pruned <- prune(rt,cp = min.xerror) |
16  fancyRpartPlot(rt.pruned)
17
```



Rattle 2016-Oct-28 18:18:44 310250154

- This tree is giving an overview of decision rules for predicting our outcome.
- We see that hour has the maximum weightage and hence it is the root of the decision/prune tree.
- Subsequent levels explain the weightage in the model respectively.

**Plots based on the best fit model:**



Residuals vs Fitted

lm(kWh ~ Temperature + month + day + `Day of Week` + hour)



Normal Q-Q

lm(kWh ~ Temperature + month + day + `Day of Week` + hour)

Residuals vs Leverage

lm(kWh ~ Temperature + month + day + `Day of Week` + hour)