

Advances in Data Science



Summary

You have been offered an internship at Tarja and Pasi Inc. in Finland and you are excited to work on your first “international” project. Tarja and Pasi Inc. is an energy modeling consultancy and they have been approached by a Vokia Inc. to help monitor and reduce energy consumption in 78 of their buildings. Vokia Inc. owns these buildings and wants to understand and reduce energy usage and wants to make the buildings more energy efficient.

Data:

You have been given two input files.

RawData.csv

BuildingID	Building	Meter number	type	date	hour	Consumption (In Kwh)
------------	----------	--------------	------	------	------	----------------------

- Note that the data covers 78 buildings and has hourly power consumption info for close to 1 year
- Use only the **elect** and **dist_heat** data which covers electricity usage and heating energy consumption data

building	address	area_floor_m.sqr(in square meter)
----------	---------	-----------------------------------

- This data has the building address (of course in Finland) and the area of the building

Table of Contents

Part 1 - Data Ingestion and Wrangling	4
1.1 - Ingest the data	
1.2 - Extract Features	
1.3 - Normalize Consumption	
1.4 - Extract the weather data	
1.5 - Merging the weather data with the ingested data	
Part 2 - Modeling Tasks	8
2.1 - Prediction	
2.1.1 - Create two calculated columns for Base hour usage & Base hour class	
2.1.2 - Build Regression Build Regression, KNN, Random Forest , Neural Network models for each building dataset (78 different models)	
2.1.3 - Build Regression, KNN, Random Forest and Neural Network models for all buildings as one dataset.	
2.1.4 - Compute MAPE, MAE, RMSE for all models. (Discuss the best model)	
2.1.5 - Outlier Detection	
2.2 - Classification - Build a model to predict the <u>Base Hour Class</u> flag for this dataset	
2.2.1 - Build Logistic Regression, KNN, Random Forest and Neural Network models for each building dataset (78 different models)	
2.2.2 - Build Logistic Regression, KNN, Random Forest and Neural Network models for all buildings as one dataset.	
2.2.3 - Compute Confusion matrices and ROC charts for all models. (Discuss)	
2.3 - Clustering	
2.3.1 - Cluster the buildings using k-means & hierarchical clustering.	

DATA INGESTION & WRANGLING

- The Data set contains 12 million entries which is hard to read as a csv file.
- So , Initially the data is read as data table and atomic variable to increase speed of processes that we do on the data.
- It is then converted to a data frame and retained only electric and dist_heating .
- Data set contains some missing data in the columns of VAC which is filled by Building names (9 & 27).
- Merge the address details along with the area per sq. unit to our data.
- To do any sort of data manipulation in a data table , we use the keyword 'setkey(datatablename,columnname)'
- Clean the data by
 - converting date column from int to date format
 - add year month and day column
 - get day of the week , weekday , peak hour

```
# Read the data set as data table and atomic variable for large dataset
mydata1<-fread("Finland_addresses_area.csv")
mydata2<-fread("Finland_masked.csv")

#convert it to dplyr data frame
dplyr::tbl_df(mydata1)
dplyr::tbl_df(mydata2)

# transpose meter number to columns
setkey(mydata1,type)
abc1<-mydata1[,c("elect", "Dist_Heating")]
Rsub1[,1:N,by=type]

#####
Rsub1<-dcast(mydata1,BuildingID+vac+date+hour+meternumb~type,value.var = "Consumption",sum)
# melt it the meter no. back to rows
Rsub1<-dplyr::melt(Rsub1[,c("BuildingID","vac","date","hour","meternumb","Dist_Heating","elect")])
#####
#assign building 9 and 27
abc1[, vac := ifelse(BuildingID==81989, "Building 27", vac)]
abc1[, vac := ifelse(BuildingID==82254|BuildingID==83427|BuildingID==84681, "Building 9", vac)]
#remove blank
abc2<-abc1[!(vac=="")
names(abc2)[2]<- "building"

#merge the address and area data with your data
setkey(abc2,building)
setkey(mydata2,building)
data_merge <- merge(abc2,mydata2, all.x=TRUE)
```

```

#convert date column from int to date format
D <- transform(Data_merge, date = as.Date(as.character(date), "%m/%d"))
#.....att11: View(Data_merge)
#add year month and day column
E<-tidyr::separate(D, date, c("year", "month", "day"))
#get date column
E$date <- as.Date(paste(E$year, E$month, E$day, sep="-")) #add date column
# get Day of week
E$DayofWeek<-as.POSIXlt(E$date)$wday
#check for weekday
E$Weekday<-ifelse(E$DayofWeek > 0 & E$DayofWeek < 6, "1", "0")
#get hour column
E$hour <- E$hour
E$hour <- as.numeric(E$hour)
E$hour <- ifelse(E$hour > 6 & E$hour < 24, "1", "0")
# convert day and month to numeric column
E$month<-as.numeric(E$month); E$day<-as.numeric(E$day)

E<-as.data.table(E)
E[,N,by=BuildingID]

```

To get list of holidays in a calendar year of 2013

- We use library(rvest) to scrape the data from the HTML page using the CSS Node
- Merge the holiday data and find the days which are holidays in our Building data .

```

#get holidays from webpage
htmlpage<-read_html("http://www.timeanddate.com/calendar/?year=2013&country=21")
#scrape CSS
holiday<-html_nodes(htmlpage,"tbl .lpad td:nth-child(1) , #tbl .col")
#get the text out of node
holiday1<-html_text(holiday)
#convert text to table
holiday2<-read.table(text = holiday1, sep = ".", colClasses = "character")
#name the column
names(holiday2)[1]<-"day"
names(holiday2)[2]<-"month"
#translate finnish to numeric month
fin_month = matrix(c(" tam", " maas", " huh", " tou", " kes", " mar", " jou", 1,3,4,5,6,11,12),nrow=7,ncol=2)
fin_month<-as.data.frame(as.matrix(fin_month))
#name the columns
names(fin_month)[1]<-"month"
names(fin_month)[2]<-"month"
fin_month$month<-as.character(fin_month$month)
fin_month$month<-as.character(fin_month$month)
#convert month to number
holidays<-merge(holiday2, fin_month, by=c("month"))
holidays$month<-1
holidays<-as.data.frame(holidays)
#make the day and month column numeric
holidays$month<-as.numeric(holidays$month); holidays$day<-as.numeric(holidays$day)
holidays<-holidays[!duplicated(holidays[,c("day", "month")]),]
#merge the holidays data and find which days are a holiday
E2<-E1
E2<- merge(E2,holidays,by=c("day","month"), all.x=TRUE,all.y = FALSE)
E2[is.na(E2)] <- 0
#final data of PART 1.1
View(E2)

```



```

#get the date range
date1 <- as.Date('2013/01/01',format = "%Y%m%d")
date2 <- as.Date('2013/12/31',format = "%Y%m%d")
# create a sequence of every day in this year
s <- seq(date, to = date2, by='days')
#modify the date format
date.replace<-str_replace_all(s,"-","/")
#get the dist number of airport
n=nrow(distinct_airport)
#scrape the html from URL
weathers<-data.frame(lime=as.character())
for(j in 1:n)
{
  code<-distinct_airport[j,1]
  dat<-paste0("https://www.ndbc.gov/history/airport/",code,"/2013/1/1/DailyHistory.html?req_city=Marikham&req_state=Finland&format=1")
  htmlpage<-read_html(dat)
  #get the text out of HTML
  weather<-html_text(htmlpage)
  #make it a table
  weath<-read.table(text = weather, sep = ",", colClasses = "character")
  #remove extra row with heading
  columns(weath)~weath[1,]
  weath = weath[-1, ]
  #name 1st column as TIME
  names(weath)[1]<-"Time"
  #convert it to data frame
  weath1<-data.frame(weath)
  #add a column of date
  weath1<- weath1 %>% rowwise() %>% mutate(date = "2013/1/1")
  #get the dist number of airport
  n=nrow(distinct_airport)
  #for loop to scrape the weather data of a particular airport code from Jan 1 2013 to Dec 31 2013
  #first loop
  for (i in 1:342)
  {
    dat<-paste0("https://www.ndbc.gov/history/airport/",code,"/date.replace[i]"/"DailyHistory.html?req_city=Marikham&req_state=Finland&format=1")
    htmlpage<-read_html(dat)
    weath<-html_text(htmlpage)
    weath<-read.table(text = weath, sep = ",", colClasses = "character")
    columns(weath)~weath[1,]
    weath = weath[-1, ]
    names(weath)[1]<-"Time"
    weath1<-data.frame(weath)
    dat = date.replace[i]
    weath1<-weath1 %>% rowwise() %>% mutate(date=dat)
    weath1<-weath1 %>% rowwise() %>% mutate(SID_CODE = CODE)

    dplyr::tbl_df(weather)
    dplyr::tbl_df(weather1)
    weathers<-dplyr::bind_rows(weathers, weath1)
  }
  weathers<-dplyr::tbl_rows(weathers, weath)
}
write_csv(weathers,"files/weathers.csv")

```


Step-wise Search Regression

```
#Forward Search
regfit.fwd=regsubsets(Normalised~.,data=trainingData, method="forward")
forwardSearchSummary=summary(regfit.fwd)
names(forwardSearchSummary)
forwardSearchSummary$rs
forwardSearchSummary$adjr2
coef(regfit.fwd,5)

par(mfrow=c(2,2))
plot(forwardSearchSummary$rs, xlab="Number of Variables ", ylab="RSS", type="l")
plot(forwardSearchSummary$adjr2, xlab="Number of Variables ", ylab="Adjusted RSq", type="l")
```

- We get similar results for stepwise regression . Therefore , we will build the linear regression model with the lowest 'p' value.

Linear Regression for the chosen independent variables

```
#Linear Regression

lm.fit = lm(Normalised ~., data = trainingData)
summary(lm.fit)
require(forecast)
accuracy(lm.fit)
L<-predict(lm.fit,testData)
accuracy(L,trainingData$Normalised)
```

Analysis of Variance Table

Response: Normalised

	DF	Sum Sq	Mean Sq	F value	Pr(>F)
Time	1	0.00	0.00	2.1571e+00	0.1419
TemperatureF	1	14.87	14.87	1.8871e+04	< 2.2e-16 ***
Dew.PointF	1	1.53	1.53	1.9414e+03	< 2.2e-16 ***
Humidity	1	0.83	0.83	1.0517e+03	< 2.2e-16 ***
Sea.level.PressureIn	1	0.00	0.00	5.2600e-02	0.8185
VisibilityMPH	1	72.10	72.10	9.1525e+04	< 2.2e-16 ***
Wind.SpeedMPH	1	2.92	2.92	3.7096e+03	< 2.2e-16 ***
Gust.SpeedMPH	1	2.27	2.27	2.8850e+03	< 2.2e-16 ***
Events	1	0.76	0.76	9.6450e+02	< 2.2e-16 ***
WindDirDegrees	1	0.73	0.73	9.2141e+02	< 2.2e-16 ***
sta_code	1	347.22	347.22	4.4077e+05	< 2.2e-16 ***
day	1	0.04	0.04	5.2783e+01	1.731e-13 ***
month	1	0.19	0.19	2.4186e+02	< 2.2e-16 ***
BuildingID	1	19.81	19.81	2.5149e+04	< 2.2e-16 ***
meternumb	1	0.23	0.23	2.8945e+02	< 2.2e-16 ***
type	1	63.05	63.05	8.0036e+04	< 2.2e-16 ***
DayofWeek	1	0.00	0.00	1.9580e-02	0.8889
Weekday	1	1.55	1.55	1.9651e+03	< 2.2e-16 ***
Peakhour	1	2.85	2.85	3.6126e+03	< 2.2e-16 ***
mont	1	0.04	0.04	4.6155e+01	1.094e-11 ***
Residuals	1	67.15	67.15	8.5737e+04	< 2.2e-16 ***
	435249	342.88	0.00		

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

> accuracy(lm.fit)

	RM	RMSE	MAE	APL	MAPE	MASE
Training set	-2.017104e-19	0.02806659	0.01402648	-67.1175	150.0256	0.6894058

- The RMSE value is 2% , which indicates the absolute fit of the model to the data and shows how accurate our dependent variable of train data set is equal to the prediction of our train data set.

KNN Regression

```
#KNN REGRESSION
sqrt(24)
trainingData1<-as.matrix(trainingData)
testData1<-as.matrix(testData)
m1<-knn(train = trainingData,test=testData,cl=train_target,k=5)
```

```
table(m1,trainingData$Normalised)
mean(test_target==m1)
```

```
summary(m1)
```

Random Forest Regression

```
#####Random forest#####

install.packages("h2o")
library(h2o)
localH2O <- h2o.init(nthreads = -1)
h2o.init()
train.h2o <- as.h2o(trainingData)
test.h2o <- as.h2o(testData)
y.dep <- 21
x.indep <- c(1:20)

system.time(
  rforest.model <- h2o.randomForest(y=y.dep, x=x.indep, training_frame =
    train.h2o, ntrees = 1000, mtries = 3, max_depth = 4, seed = 1122)
)

td1_normalise<-as.data.frame(testData[[21]])
x2<-data.frame(sub_rf,td1_normalise)
plot(x2)

#check variable importance
h2o.varimp(rforest.model)

#making predictions on unseen data
system.time(predict.rforest <- as.data.frame(h2o.predict(rforest.model, test.h2o)))
sub_rf <- data.frame(BaseHourRate = predict.rforest$predict)
```

- Running random forest model in R threw an error a “memory heap error”
- To optimize it , we used the library(h2o) which uses the JVM to execute the code ,hence removing the memory issue.

Variable Importances:

	variable	relative_importance	scaled_importance	percentage
1	sta_code	139496.281250	1.000000	0.450595
2	BuildingID	95717.132812	0.686163	0.389182
3	Gust.SpeedNPH	27245.988281	0.195317	0.088009
4	VisibilityNPH	25553.988281	0.183188	0.082543
5	TemperatureF	3666.728027	0.026285	0.011844
6	meternumb	3629.106934	0.026016	0.011723
7	Dew.PointF	3580.954834	0.025097	0.011309
8	type	3336.192871	0.023916	0.010776
9	month	2840.497803	0.020363	0.009175
10	Events	1261.291748	0.009042	0.004074
11	Wind.SpeedNPH	874.869751	0.006272	0.002826
12	Sea.Level.PressureIn	768.606689	0.005510	0.002483
13	DayofWeek	332.196564	0.002381	0.001073
14	Weekday	321.923981	0.002308	0.001040
15	Peakhour	318.687805	0.002285	0.001029
16	Humidity	251.027344	0.001800	0.000811
17	WindDirDegrees	202.399017	0.001451	0.000654
18	Time	112.111389	0.000804	0.000362
19	day	91.192726	0.000654	0.000295
20	mont	61.033161	0.000438	0.000197

- From the linear regression model , we know that we need 5 variables for the best results. So we consider the top 5 variables from the above figure for our model.
- It is constant with the variables deemed fit by linear regression model.

Conclusion for Prediction

For regression model , linear regression using `lm.fit` and random forest will give accurate models as compared to k nearest neighbor . The neural network as seen below can be difficult to interpret but with forecast it might yield better results.

Classification

Exhaustive Search Regression

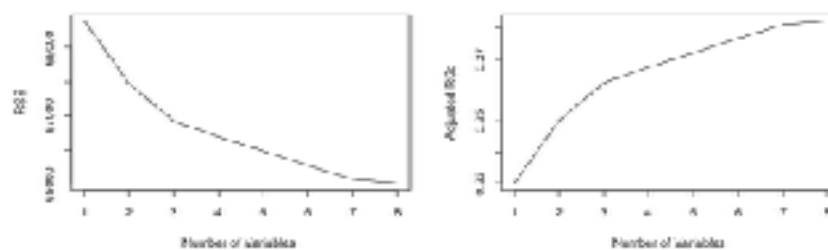
```
#Regression
#Exhaustive Search
regfit.full=regsubsets(BaseHourRate~,data=trainingData,really.big = TRUE)
exhaustiveRegSummary =summary(regfit.full)
print(exhaustiveRegSummary)
par(mfrow=c(2,2))
plot(exhaustiveRegSummary$rsr, xlab="Number of Variables", ylab="RSS", type="l")
plot(exhaustiveRegSummary$adjr2, xlab="Number of Variables", ylab="Adjusted RSq", type="l")
```

```
Selection Algorithm: exhaustive
Time Temperature Dew point Humidity Sea Level Pressure Visibility Wind Speed Wind Dir. SpeedDir
1 [ 3 ] 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0
2 [ 3 ] 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0
3 [ 3 ] 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0
4 [ 3 ] 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0
5 [ 3 ] 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0
6 [ 3 ] 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0
7 [ 3 ] 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0
8 [ 3 ] 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0

Dew pt WindDir Dew pt WindDir Dew pt WindDir Dew pt WindDir Dew pt WindDir Dew pt WindDir Dew pt WindDir
1 [ 3 ] 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0
2 [ 3 ] 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0
3 [ 3 ] 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0
4 [ 3 ] 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0
5 [ 3 ] 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0
6 [ 3 ] 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0
7 [ 3 ] 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0
8 [ 3 ] 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0

Some Load
1 [ 3 ] 0.0
2 [ 3 ] 0.0
3 [ 3 ] 0.0
4 [ 3 ] 0.0
5 [ 3 ] 0.0
6 [ 3 ] 0.0
7 [ 3 ] 0.0
8 [ 3 ] 0.0
```

- Looking at the selection algorithm we can come to a conclusion that 8 independent variables have major influence on the Base Hour Class (Variable name : BaseHourRate)
- The below image confirms our conclusion.



- Forward Search Classification rendered the same result.

Logistic Regression

#Logistic Regression

```
glm.fit = glm(BaseHourRate ~., data = trainingData,family = "binomial")
summary(glm.fit)
accuracy(glm.fit)
L<-predict(glm.fit,testData)
accuracy(L,trainingData$Normalised)
```

```
Coefficients:
            Estimate Std. Error z value Pr(>|z|)
(Intercept)  2.767e-05  9.968e+04  2.017 8.24e-05 ***
Time         2.438e-03  1.543e+04  6.118  0.875
TemperatureF -1.679e-04  1.077e+05  -6.106  0.876
Dew PointF   1.609e-04  1.140e+05  6.140  0.880
Humidity     -4.095e-03  4.288e+04  -6.096  0.924
Sea Level PressureIn -9.840e-04  6.177e+04  -1.544  0.123
VisibilityMPH -3.345e-03  7.513e+04  -6.045  0.964
Wind.SpeedMPH 7.772e-03  3.129e+04  6.218  0.804
Gust.SpeedMPH -2.575e-03  2.720e+04  -6.095  0.925
Events       -3.127e-03  6.752e+04  -6.046  0.963
WindDirDegrees -2.477e-02  1.688e+03  -6.228  0.820
sta_code     -1.187e-09  2.535e+08  -4.045 5.23e-05 ***
day          -2.741e-03  1.181e+04  -6.232  0.816
month        6.150e-03  3.589e+04  6.172  0.804
Building_U   -1.450e-01  1.223e+01  -1.204  0.225
meternum     1.215e-07  1.009e+07  6.711  0.477
type         -7.330e-03  6.204e+05  -6.800  0.373
DaysOffWeek  -4.118e-03  3.831e+04  -6.078  0.937
Weekday      -2.027e-04  2.720e+05  -6.104  0.917
Poolhour     2.890e-04  2.618e+05  6.110  0.912
month        1.314e-05  1.510e+05  6.154  0.877
Normalised   6.700e-05  1.134e+05  5.452 5.10e-01 ***
```

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

```
(Dispersion parameter for binomial family taken to be 1)
```

```
Null deviance: 5.2843e+05 on 435270 degrees of freedom
Residual deviance: 6.1757e-01 on 435249 degrees of freedom
AIC: 46.418
```

```
Number of Fisher Scoring iterations: 25
```

- Looking at the 'p' values of the logistic regression , we can infer that 'sta_code' & 'Normalised' variables influence the dependent variable (Base Hour Rate) highly.


```

> accuracy(glm.fit)
              ME          RMSE          MAE  MPE MAPE          MASE
Training set -1.856771e-08 5.386075e-05 7.079525e-07 -Inf  Inf 1.700283e-06
> L<-predict(lm.fit,testData)
> L<-predict(glm.fit,testData)
> accuracy(L,trainingData$Normalised)
              ME          RMSE          MAE  MPE  MAPE
Test set 439.3712 29042.76 13095.64 5945148 162620900
> |

```

- From the predictions RMSE , we have come to the following conclusion
 - our division of test and train data with respect to quantity might be wrong.
 - looking at the sta_code , we have observed that there are two distinct station codes(airport code) and one airport code is very small with respect to the other . This could possibly influence the prediction as the division of data might not be even.
 - use of multiple numeric independent variable to obtain regression of factor dependent variable might have resulted in bad model.

KNN Classification

```
#KNN Classification
sqrt(24)
m2<-knn(train = trainingData,test=testData,cl=train_target,k=5)
table(test_target,m2)
```

```
> table(test_target,m2)
      m2
test_target  0      1
0 131664      0
1      0 54881
```

```
> mean(test_target==m2)
[1] 1
```

- it looks like the algorithm successfully predicted base hour rate for our test data with an accuracy of 100%.
- this results could be because we have normalized our input variable between a constant range of values (0 - 1)

Random Forest

```
#Random Forest classification

y.depl <- 22
x.indepl <- c(1:21)

system.time(
  rforest.modell <- h2o.randomForest(y=y.depl, x=x.indepl, training_frame = train.h2o,
                                     ntree = 1000, mtries = 3, max_depth = 4, seed = 1122)
)
#check variable importance
h2o.varimp(rforest.modell)

#making predictions on unseen data
system.time(predict.rforest1 <- as.data.frame(h2o.predict(rforest.modell, test.h2o)))
sub_rf <- data.frame(BaseHourRate = predict.rforest1$predict)
```

Variable Importances:

	variable	relative_importance	scaled_importance	percentage
1	Normalised	29537094.000000	1.000000	0.840784
2	sta_code	1431060.125000	0.048450	0.040736
3	BuildingID	928483.750000	0.031434	0.026430
4	month	694171.687500	0.023502	0.019750
5	TemperatureF	528416.625000	0.017890	0.015042

	variable	relative_importance	scaled_importance	percentage
16	Sea.Level.PressureIn	14193.941406	0.000481	0.000404
17	Humidity	9323.667969	0.000316	0.000265
18	Wind.SpeedMPH	7883.343262	0.000267	0.000224
19	day	2361.557617	0.000080	0.000067
20	WindDirDegrees	1974.001343	0.000067	0.000056
21	mont	1511.108887	0.000051	0.000043

>

	BaseHourRate	testDataA..22..
1	0	0
2	0	0
3	1	1
4	1	1
5	0	0
6	1	1

- Look at our predictions for random forest for classification , we are getting 100% accuracy for predicting the test data with train data for Base Hour rate

Conclusion for Classification :

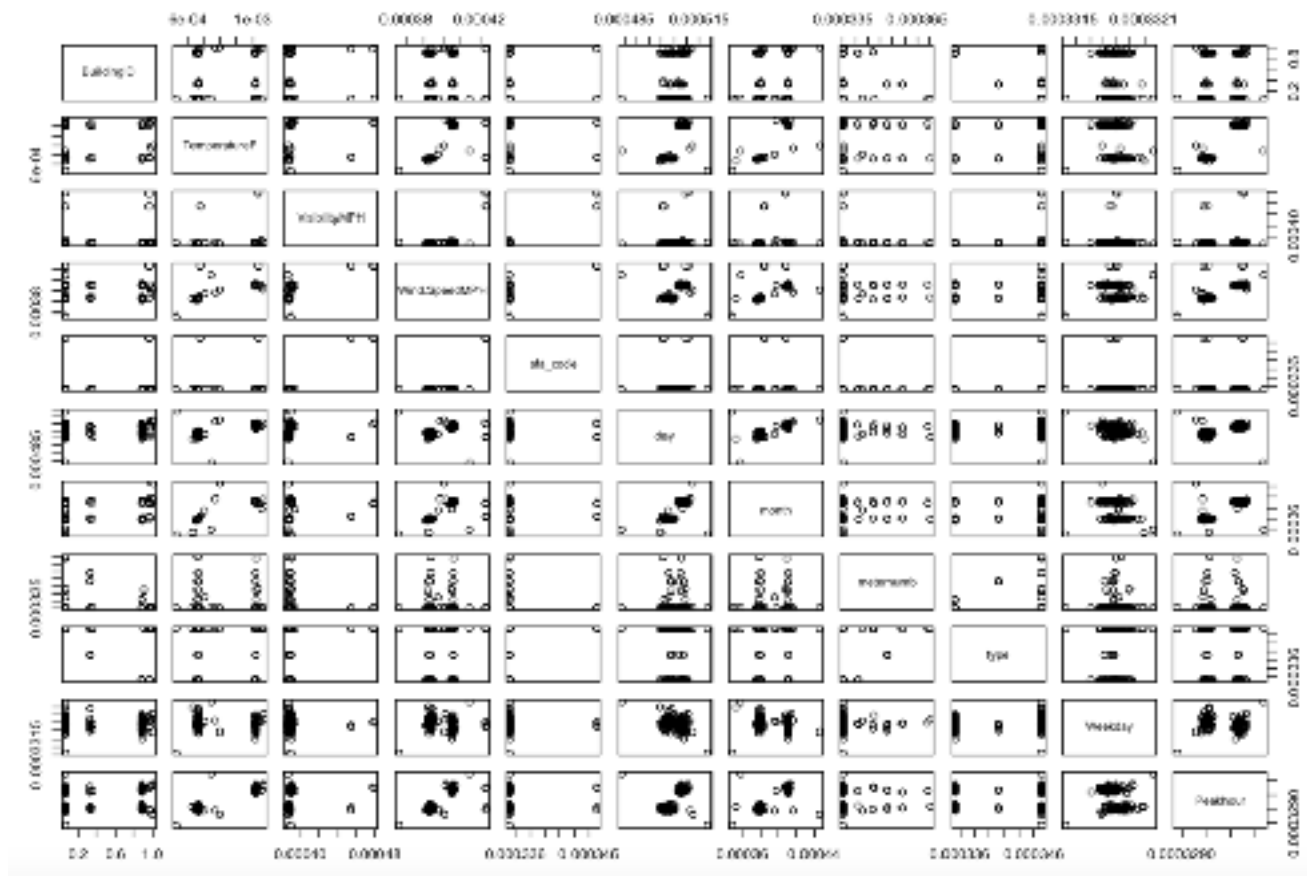
For the regression model , K nearest neighbor and Random Forest yielded high accuracy as compared to logistic model.

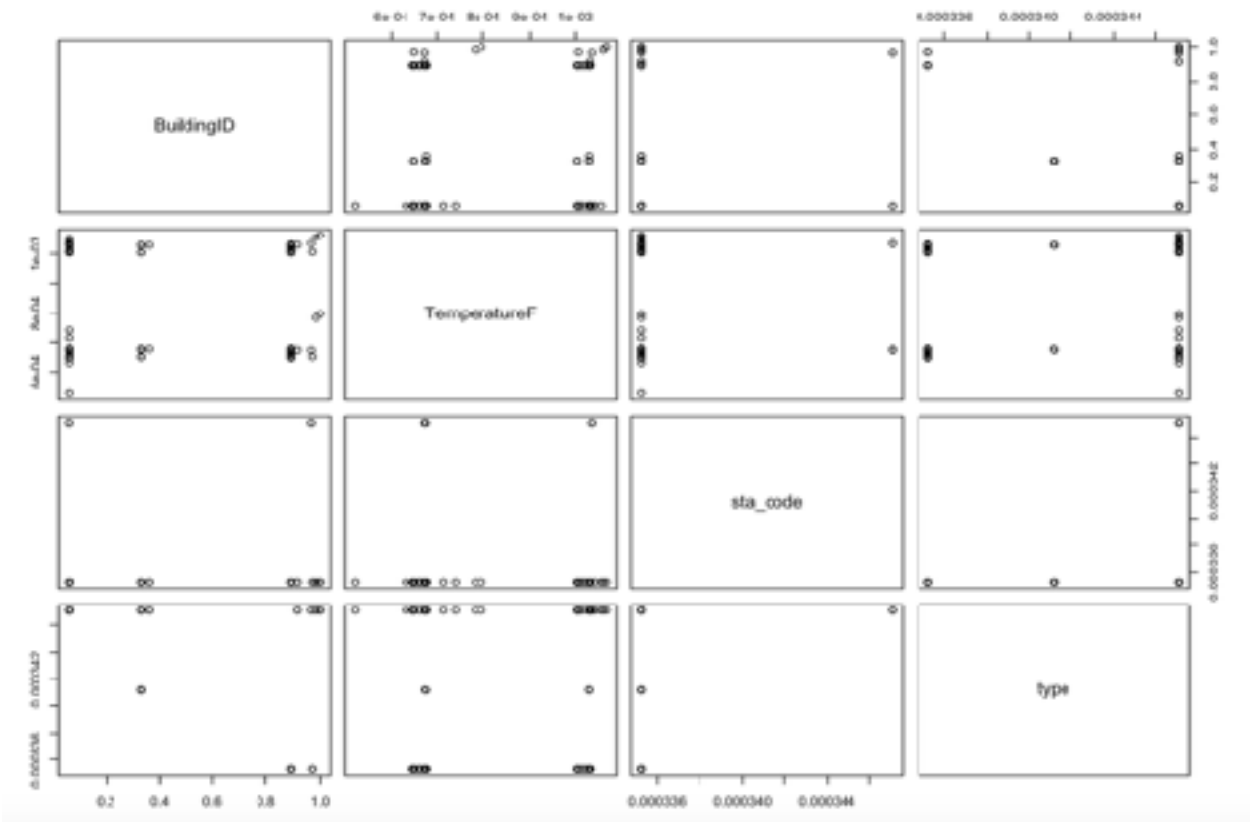
Clustering

K-means Clustering

```
#Kmeans clustering
fx <- trainingData[,data.table(kmeans(cbind(TemperatureF,sta_code,type),centers=2)$centers),by=BuildingID]
fx1 <- trainingData[,data.table(kmeans(cbind(TemperatureF,VisibilityMPH,Wind.SpeedMPH,sta_code,day,
month,meternumb,type,Weekday,Peakhour),
centers=2)$centers),by=BuildingID]

plot(fx1)
plot(fx)
```

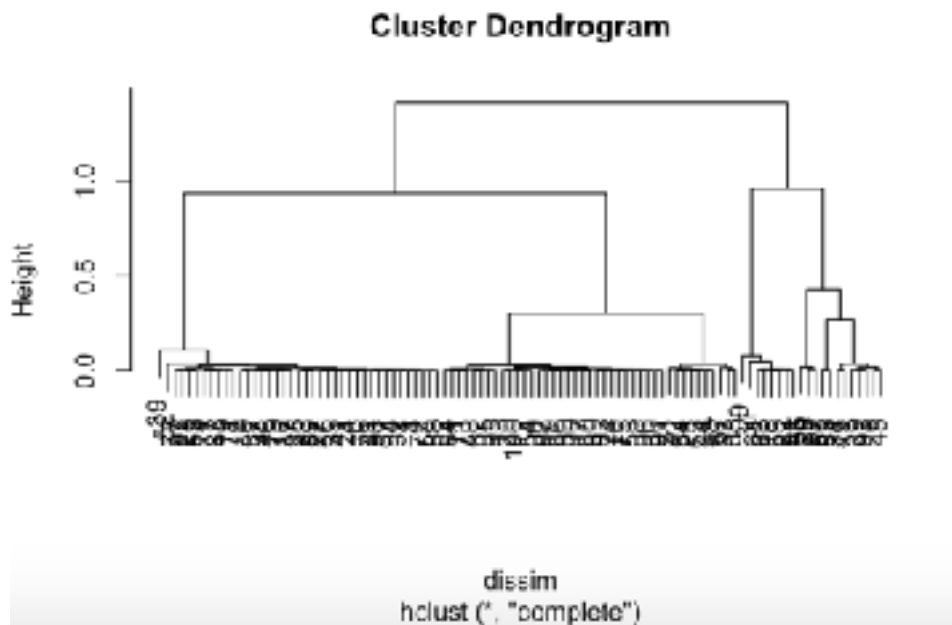




Hierarchical Clustering

```
#Hierarchial clustering
install.packages("fastcluster")
library('fastcluster')
df<-as.data.frame(trainingData)
dissim <- dist(df[1:50,1:10])
hr <- hclust(dissim)
plot(hr)

clusterCut <- cutree(hr, 4)
table(clusterCut, df$BuildingID[1:50])
```



NEURAL NETWORK

Artificial Neural Networks(ANN) are implemented by the following process,

1. Do the train-test split and fit the model to the train set
2. Normalize the data
3. Categorize the input parameters based on regression
4. Implement the neuralnet() function
5. Carry out multiple iterations.

1.After getting the input from part 1, we start by randomly splitting the data into a train and a test set, then we fit a linear regression model and test it on the test set.

```
index <- sample(1:nrow(DT_2),round(0.75*nrow(DT_2)))
train <- DT_2[index,]
test <- DT_2[-index,]
lm.fit <- glm(Consumption~., data=train)
summary(lm.fit)
pr.lm <- predict(lm.fit,test)
MSE.lm <- sum((pr.lm - test$Consumption)^2)/nrow(test)
```

2. Normalizing the data using the min-max method.

```
maxs <- apply(DT_2, 2, max)
mins <- apply(DT_2, 2, min)
scaled <- as.data.frame(scale(DT_2, center = mins, scale = maxs - mins))
train_ <- scaled[index,]
test_ <- scaled[-index,]
```

3.Input parameters based on the regression

- Weekday
- Peakhour
- Mont
- Sta_code

- Day
- Month
- BuildingID
- Meternumb
- Type
- Area_floor_m.sqr
- TemperatureF
- Wind.SpeedMPH
- Gust.SpeedMPH

4 Implementing the function `neuralnet()` for the train and test data.

```
nn <- neuralnet(
  Consumption ~Weekday+Peakhour+mont+sta_code+day+month+
    BuildingID+meternumb+type+TemperatureF+Wind.SpeedMPH+Gust.SpeedMPH,
  data=train_, hidden=5, err.fct="ce",
  linear.output=FALSE)

nn <- neuralnet(
  Consumption ~Weekday+Peakhour+mont+sta_code+
    day+month+BuildingID+meternumb+type+TemperatureF+Wind.SpeedMPH+Gust.SpeedMPH,
  data=test_, hidden=5, err.fct="ce",
  linear.output=FALSE)
```

5a. When we used one neuron in the hidden layer, we found the error to be very high, 318 to be precise.

The result and the pictorial representation is as shown below.

```
error          318.643978590810
reached.threshold 0.009923574821
steps          11343.000000000000
Intercept.to.1layhid1 -16.199244005901
Weekday.to.1layhid1 -0.797269070692
Peakhour.to.1layhid1 -1.035939020437
mont.to.1layhid1 0.227693918439
sta_code.to.1layhid1 38.499532668095
day.to.1layhid1 -0.237598592023
month.to.1layhid1 0.762361266469
BuildingID.to.1layhid1 13.556366162686
meternumb.to.1layhid1 16.813724278667
type.to.1layhid1 15.446192809758
TemperatureF.to.1layhid1 10.469788211412
Wind.SpeedMPH.to.1layhid1 -1.549429791468
```

```
Gust.SpeedMPH.to.1layhid1  0.727833782168
Intercept.to.Consumption    -1.560359741601
1layhid.1.to.Consumption    -1.897422456452
```

5b. We used 5 neurons in the hidden layer and we find the error to be 1.31, a very significant improvement from our previous iteration. Our model now is way more effective and precise.

```
error                        1.318237568931
reached.threshold           0.004516627377
steps                       19.000000000000
Intercept.to.1layhid1       2.855457390283
Weekday.to.1layhid1         -0.212721260154
Peakhour.to.1layhid1        0.363186786673
mont.to.1layhid1            -0.956285792057
sta_code.to.1layhid1        0.365870527727
day.to.1layhid1             -1.369086711318
month.to.1layhid1           -0.720694728398
BuildingID.to.1layhid1      2.642342343749
meternumb.to.1layhid1       -0.536251391682
type.to.1layhid1            -0.928107850618
TemperatureF.to.1layhid1    -1.469454309811
Wind.SpeedMPH.to.1layhid1   -1.561321086177
Gust.SpeedMPH.to.1layhid1   0.230634231519
Intercept.to.1layhid2       1.477126460176
Weekday.to.1layhid2         0.904356922957
Peakhour.to.1layhid2        3.031855605898
mont.to.1layhid2            2.583836778537
sta_code.to.1layhid2        0.683971789602
day.to.1layhid2             -0.477954798895
month.to.1layhid2           -0.581897516734
BuildingID.to.1layhid2      -0.820341771254
meternumb.to.1layhid2       0.725141469330
type.to.1layhid2            0.991061798303
TemperatureF.to.1layhid2    0.905089119727
Wind.SpeedMPH.to.1layhid2   2.053847914534
Gust.SpeedMPH.to.1layhid2   -0.170235618395
Intercept.to.1layhid3       1.977353993978
Weekday.to.1layhid3         0.732303865565
Peakhour.to.1layhid3        0.777574029014
mont.to.1layhid3            -0.803784782927
sta_code.to.1layhid3        -0.669269731634
day.to.1layhid3             -0.663888073167
month.to.1layhid3           -0.206355884779
BuildingID.to.1layhid3      1.111310851212
```

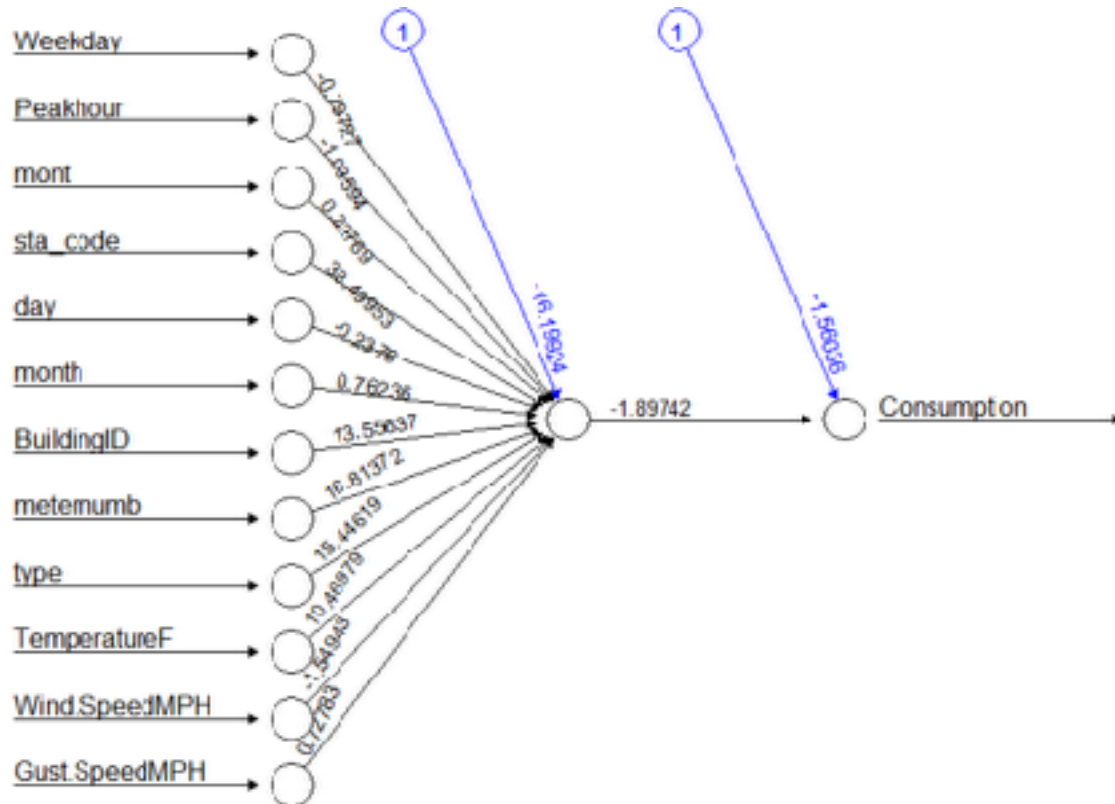
```

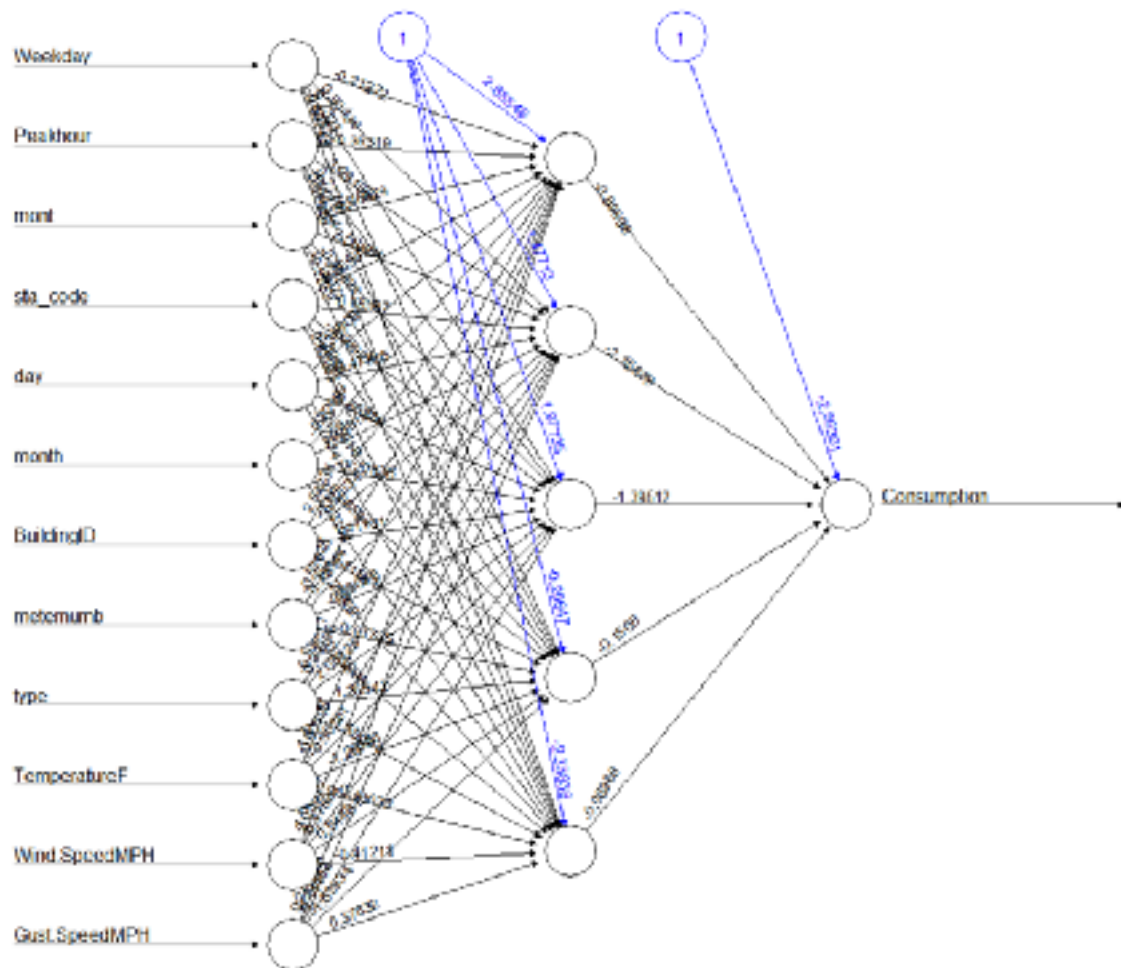
meternumb.to.1layhid3  -0.050778871185
type.to.1layhid3      2.051254070581
TemperatureF.to.1layhid3 -0.003813760523
Wind.SpeedMPH.to.1layhid3 0.378222389006
Gust.SpeedMPH.to.1layhid3 0.144818069673
Intercept.to.1layhid4  -0.296466875766
Weekday.to.1layhid4    0.448246274196
Peakhour.to.1layhid4  -0.443340483840
mont.to.1layhid4       0.190705890739
sta_code.to.1layhid4   -0.825902717874
day.to.1layhid4        -0.355175507222
month.to.1layhid4      -0.085313400164
BuildingID.to.1layhid4 0.415807489276
meternumb.to.1layhid4  -0.617353197524
type.to.1layhid4       1.355422512525
TemperatureF.to.1layhid4 -1.298482875283
Wind.SpeedMPH.to.1layhid4 0.544811243808
Gust.SpeedMPH.to.1layhid4 -1.639336655807
Intercept.to.1layhid5  -2.338363117352
Weekday.to.1layhid5    0.537142876512
Peakhour.to.1layhid5   0.819459306233
mont.to.1layhid5       -0.395904640687
sta_code.to.1layhid5   0.819848519567
day.to.1layhid5        -2.703404111654
month.to.1layhid5      -1.637540210180
BuildingID.to.1layhid5 -0.268324558076
meternumb.to.1layhid5  0.959948101243
type.to.1layhid5       -1.550136299243
TemperatureF.to.1layhid5 -0.490328732288
Wind.SpeedMPH.to.1layhid5 -0.412110686432
Gust.SpeedMPH.to.1layhid5 0.378315447403
Intercept.to.Consumption -2.563012910893
1layhid.1.to.Consumption -0.895376470240
1layhid.2.to.Consumption -2.156791886898
1layhid.3.to.Consumption -1.395115463591
1layhid.4.to.Consumption -0.150803398059
1layhid.5.to.Consumption -0.069682441638

```

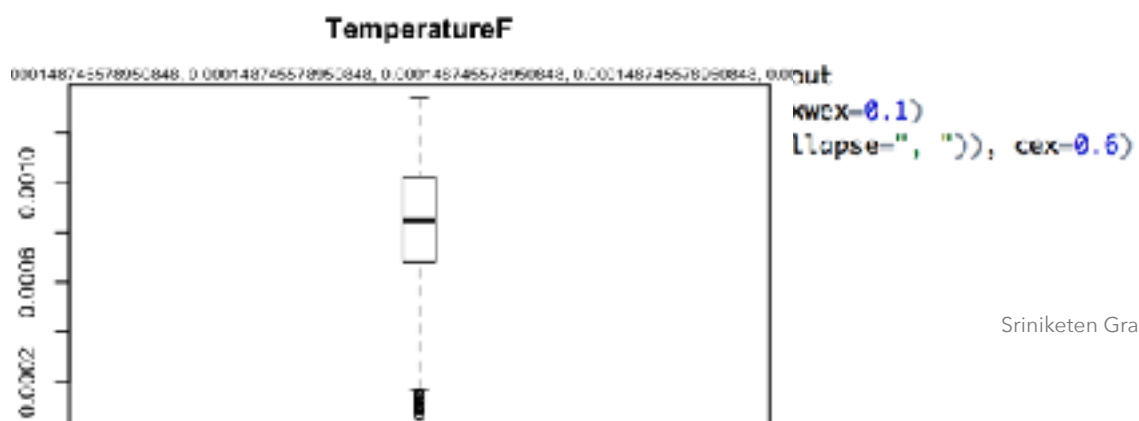
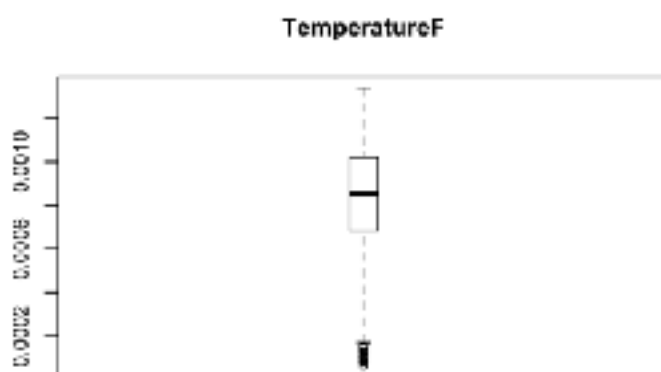
The result above displays the weights of the paths from the input to the hidden layer and from the hidden layer to the output layer. The plot below gives us a pictorial representation of the result.

The result is plotted by using the script, plot(nn)





OUTLIER DETECTION



The outliers which are not too unrealistic are retained even though they are out of the box plot even though they are quite a lot in number.

On visual examination , we have removed -9999F as temperature value can never be in that value range.