

QIANG SUN/21804416

1. Explore some obvious rules and patterns from the data.
2. Build a model to predict power consumption for each month.
3. Build a model to figure out the influence on power factor from different sub meters.
4. Find rules in the power consumption.

1. Data frame pre-process
  - a. Load it as data frame with R language
  - b. Assign each column with proper type
  - c. Assign a column to control the version and data provenance

First, we use summary to check the portion of the missing data for each column.

Fig 1: Summary of raw data

Then we examine the whole data set to check the distribution of missing data and the whole data set. We need to find whether the missing data is random or not, and find a solution for the missing data. We can see it from the Fig 2

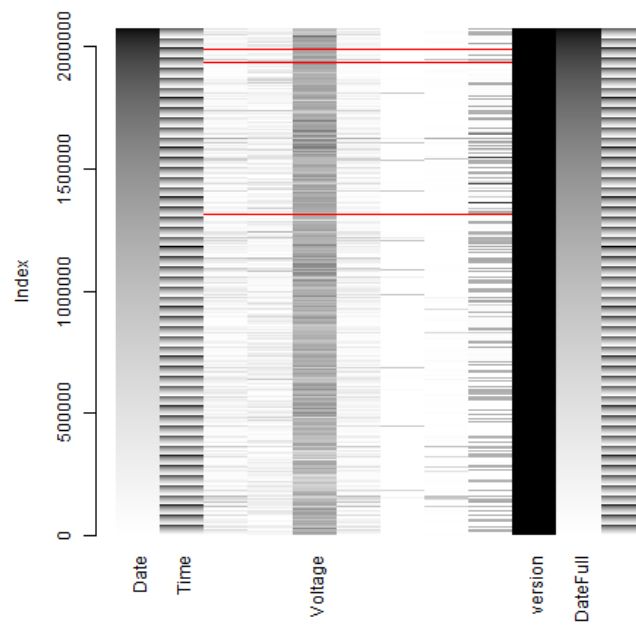


Fig 2.1 The missing data and the distribution of all the data.

The missing data is in red, the colour depth represents the value of the data.

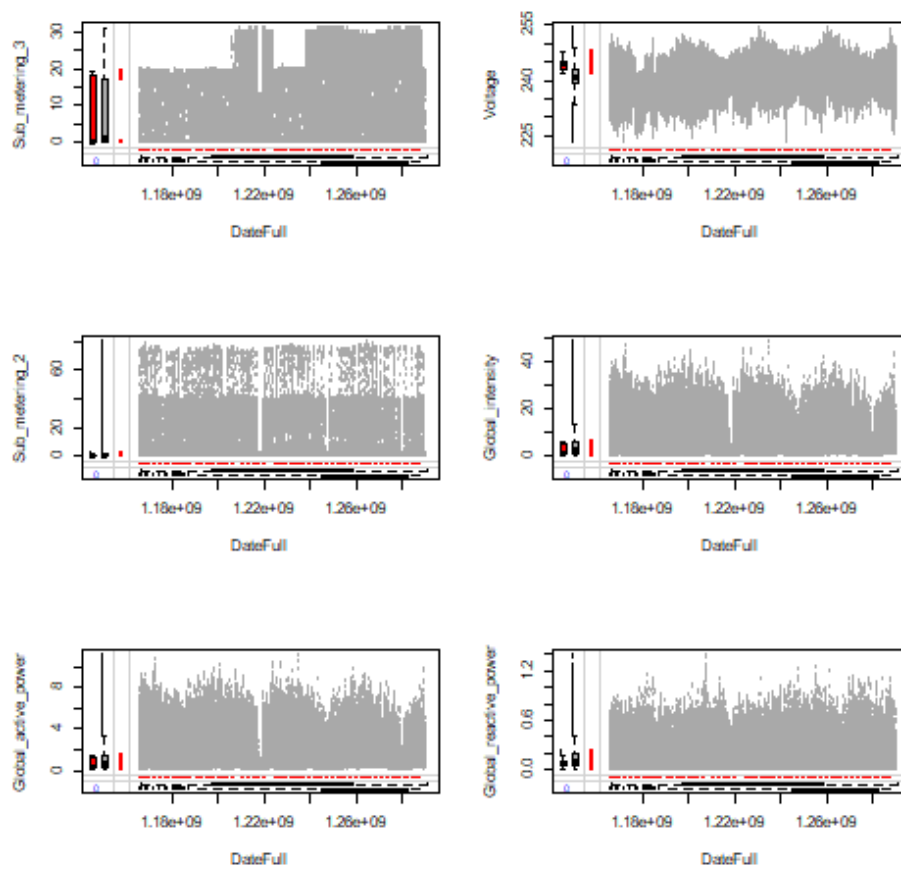


Fig 2.2 The distribution of Missing data

We can find that the missing data of time are mostly laying back of the whole date, which most in 2007 and 2008. And the distribution of the other data is mostly the same distribution as the original ones.

After analysing all the missing data, we can infer that the missing data may be caused by massive outage. And it will not affect the whole distribution of the data set if we drop them all. And we have no way but dropping all the missing data.

## Explore Data

### 1. Data Process

- Add the column: power consumption which not recorded by the 3 meters.
- Separate the date for Year, Month, Day, Hour, Minutes, and Day of Year
- Calculate apparent power, and record it in Global\_Power
- Calculate power factor, and record it in Phase

### 2. Explore the data range

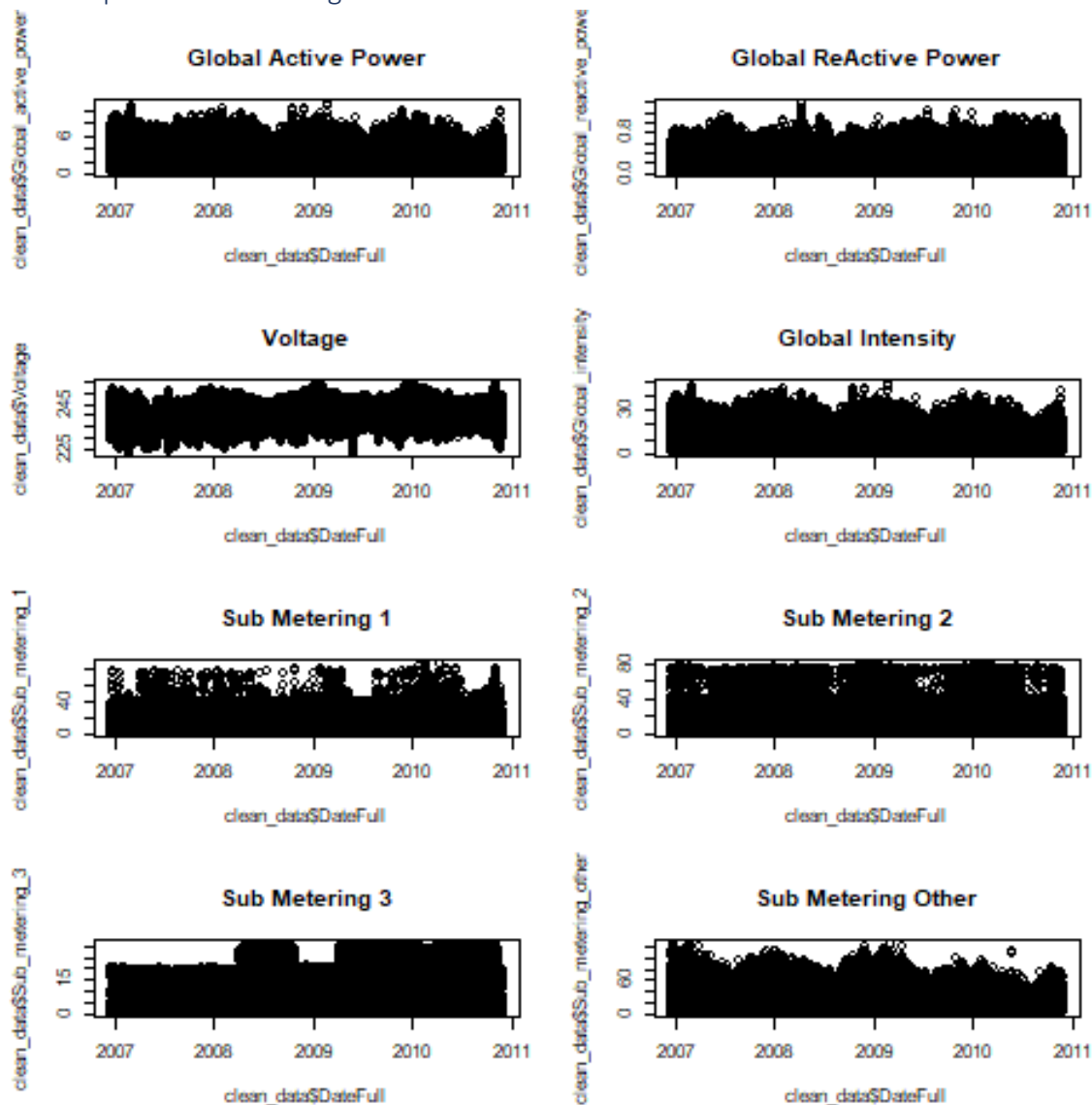


Fig 3.1 Data Range for all the data columns

We can infer from the Fig 3.1 that, the power consumption should be the lowest in June, and highest in December. Maybe the household was out for holiday in every June. And we can infer that the household is living in Southern Hemisphere. Voltage is distributed between 210-260. And there are some obvious patterns from the 3 meters. The summary of cleaned data supports our inferences.

```
> summary(clean_data)
      Date      Time      Global_active_power Global_reactive_power Voltage Global_intensity Sub_metering_1 Sub_metering_2
Min. :2006-12-16 19:51:00 1427      Min. : 0.076      Min. :0.0000      Min. :223.2      Min. : 0.200      Min. : 0.000      Min. : 0.000
1st Qu.:2007-12-10 19:52:00 1427      1st Qu.: 0.308      1st Qu.:0.0480      1st Qu.:239.0      1st Qu.: 1.400      1st Qu.: 0.000      1st Qu.: 0.000
Median :2008-11-30 19:53:00 1427      Median : 0.602      Median :0.1000      Median :241.0      Median : 2.600      Median : 0.000      Median : 0.000
Mean :2008-12-01 19:54:00 1427      Mean : 1.092      Mean :0.1237      Mean :240.8      Mean : 4.628      Mean : 1.122      Mean : 1.299
3rd Qu.:2009-11-23 19:55:00 1427      3rd Qu.: 1.528      3rd Qu.:0.1940      3rd Qu.:242.9      3rd Qu.: 6.400      3rd Qu.: 0.000      3rd Qu.: 1.000
Max. :2010-11-26 19:56:00 1427      Max. :11.122      Max. :1.3900      Max. :254.2      Max. :48.400      Max. :88.000      Max. :80.000
      (other) :2040598
Sub_metering_3 version DateFull DateTimes Sub_metering_other Year Month
Min. : 0.000      Min. :3      Min. :2006-12-16 17:24:00      Min. :2017-10-01 00:00:00      Min. : 0.000      Length:2049160      Length:2049160
1st Qu.: 0.000      1st Qu.:3      1st Qu.:2007-12-10 06:07:45      1st Qu.:2017-10-01 06:00:00      1st Qu.: 3.800      Class :character      Class :character
Median : 1.000      Median :3      Median :2008-11-30 02:22:30      Median :2017-10-01 12:00:00      Median : 5.500      Mode :character      Mode :character
Mean : 6.459      Mean :3      Mean :2008-12-02 02:01:04      Mean :2017-10-01 11:59:46      Mean : 9.315
3rd Qu.:17.000      3rd Qu.:3      3rd Qu.:2009-11-23 21:01:15      3rd Qu.:2017-10-01 18:00:00      3rd Qu.:10.367
Max. :31.000      Max. :3      Max. :2010-11-26 21:02:00      Max. :2017-10-01 23:59:00      Max. :124.833

      Day      Hour      Minute      DayOfYear      Global_Power      Phase
Length:2049160      Length:2049160      Length:2049160      Length:2049160      Min. : 0.0760      Min. :0.5559
Class :character      Class :character      Class :character      Class :character      1st Qu.: 0.3319      1st Qu.:0.9520
Mode :character      Mode :character      Mode :character      Mode :character      Median : 0.6340      Median :0.9934
                        Mean : 1.1096      Mean :0.9637
                        3rd Qu.: 1.5385      3rd Qu.:0.9997
                        Max. :11.1234      Max. :1.0000
```

Fig 3.2 Summary of Data which have been cleaned.

### 3. Some Interesting Patterns

#### a. Voltage and Intensity VS Power Factor

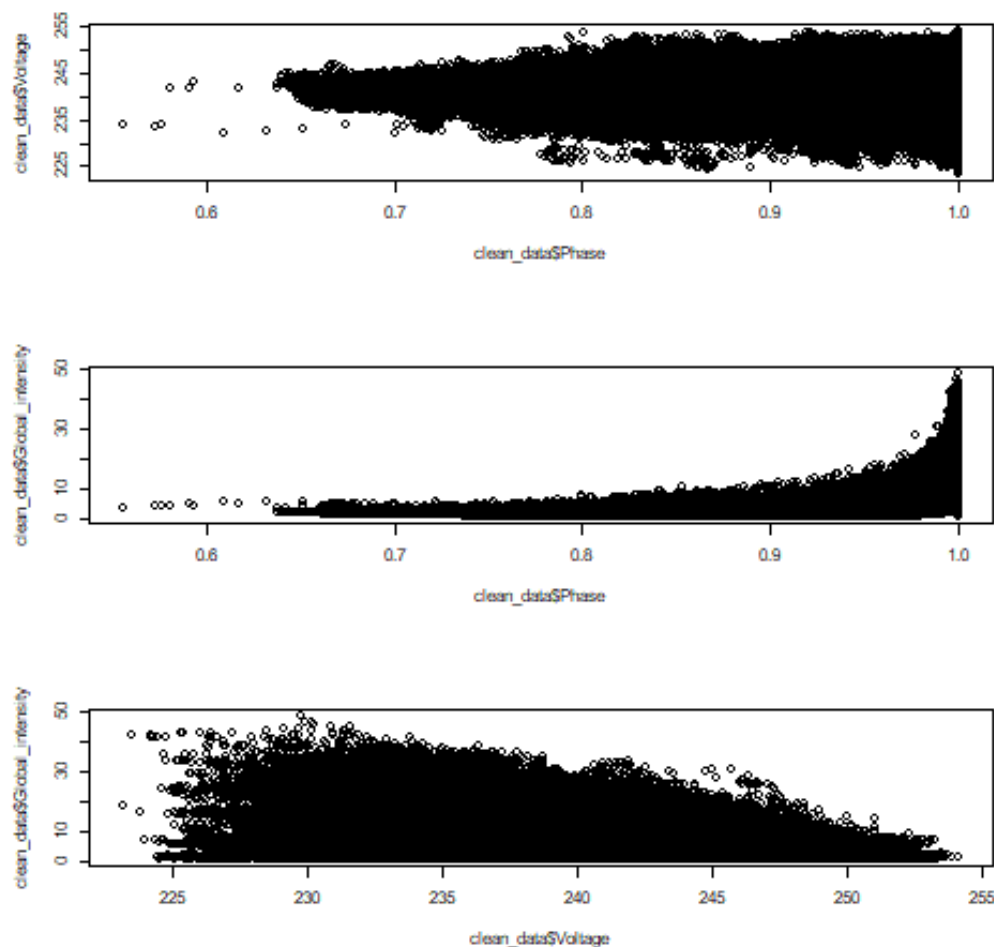


Fig 4.1 Some patterns between voltage, intensity and power factor

The distribution for Power Factor and Voltage is normal distribution, and the distribution for Power Factor and Intensity is exponential distribution. These 2 patterns can be useful for power providers to provide stable voltage and better electric quality, or used to prevent intensity to be too high, which could cause fire or grid crash.

*b. Proportion of different meters*

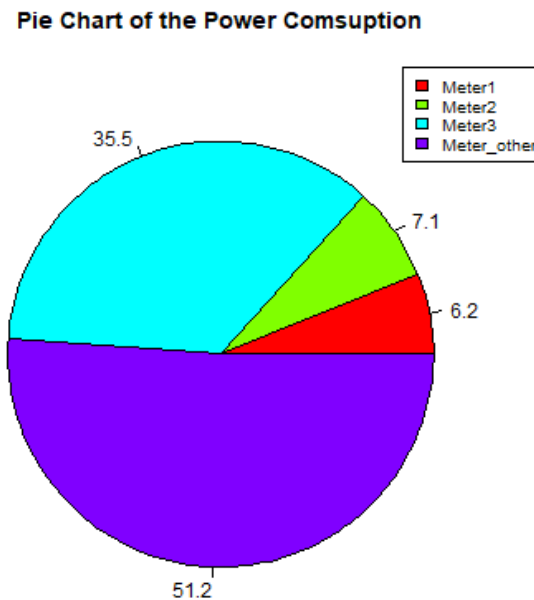


Fig 4.2 Proportion of different metters

Most power are consumed by the other equipment, and around half of the power are consumed by the 3 meters, Air Condition and Heater consume most of the power.

If the household want to save their cost of electric, they should decrease the usage of AC and heater, or they should replace a new equipment which consume less power.

## Build Models

### 1. Regression: No Linear

We want to predict the power consumption for days or months, so we can predict the cost of electric for the household in the future.

To fulfil that, we first group the power consumption data with day and month. And do the boxplots to check which version is better to fitting.

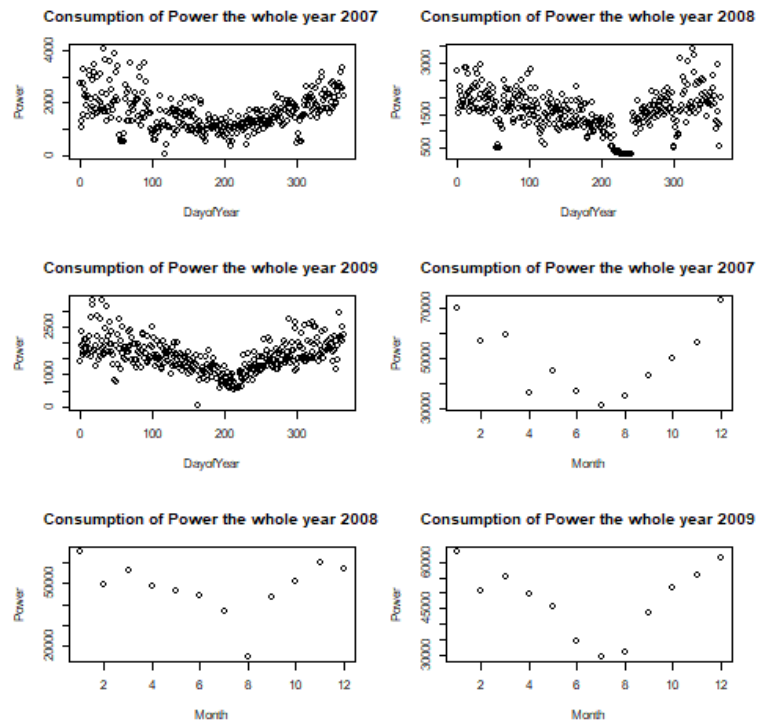


Fig 5.1 Scatter Plots of Month and Day in 2007,2008,2009

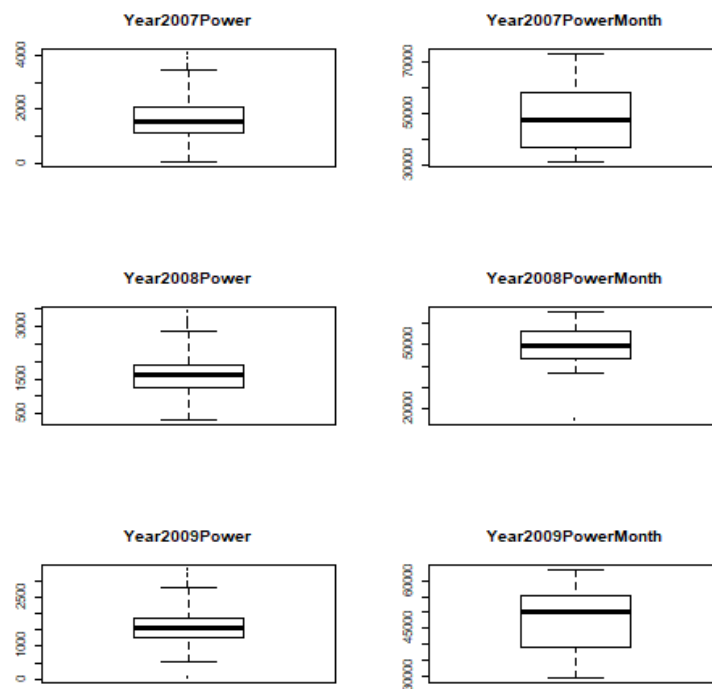


Fig 5.2 Boxplots for Power Consumption

After examining the boxplots and scatter plots, we decide to fit month vs Power Consumption with no linear regression.

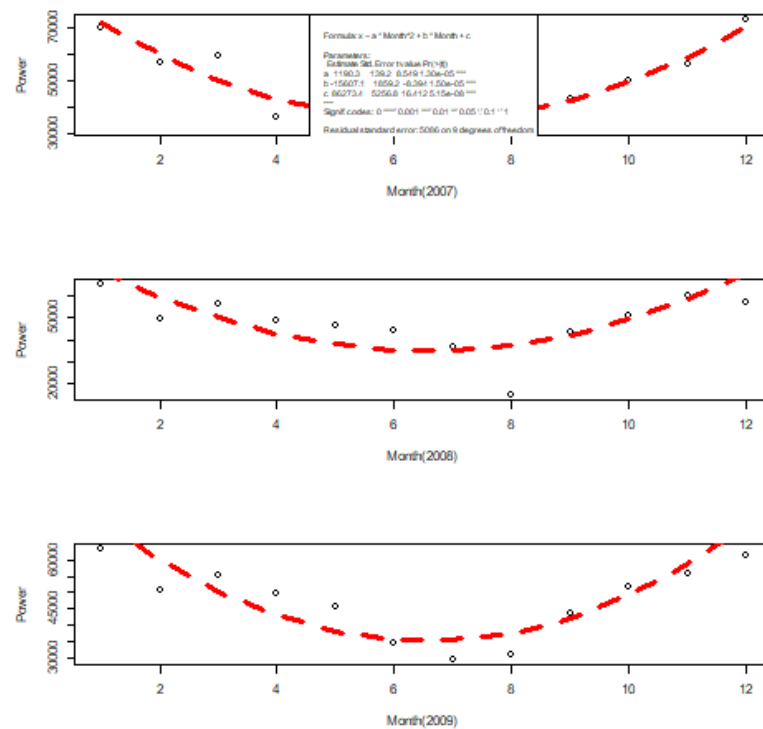


Fig 5.3 Fitting Plot

It fits pretty good. The equation is:  $\text{Power} = 1190.3 \times \text{Month}^2 - 156007.1 * \text{Month} + 86273.4$

## 2. Regression: Multivariable Linear

The higher the factor power is, the less power waste, so we want to find out which sub meter influence the power factor most. With that power provider can

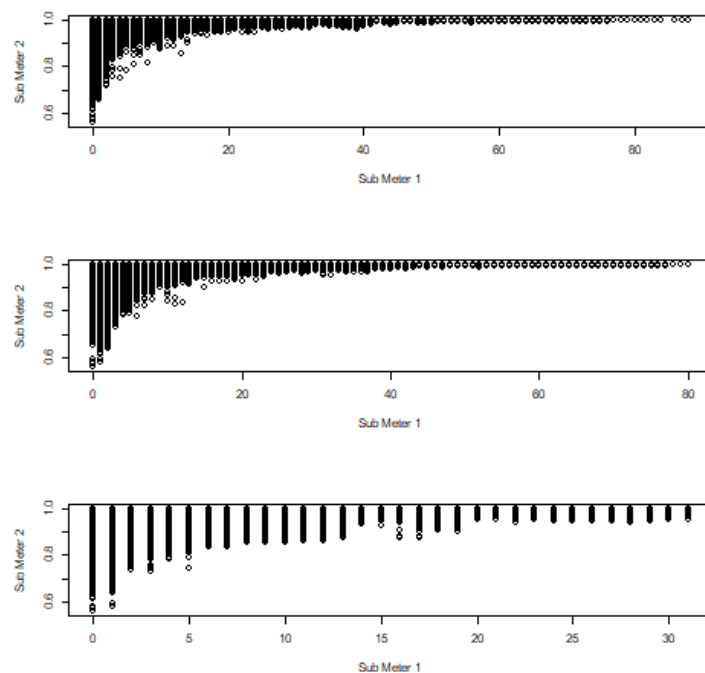
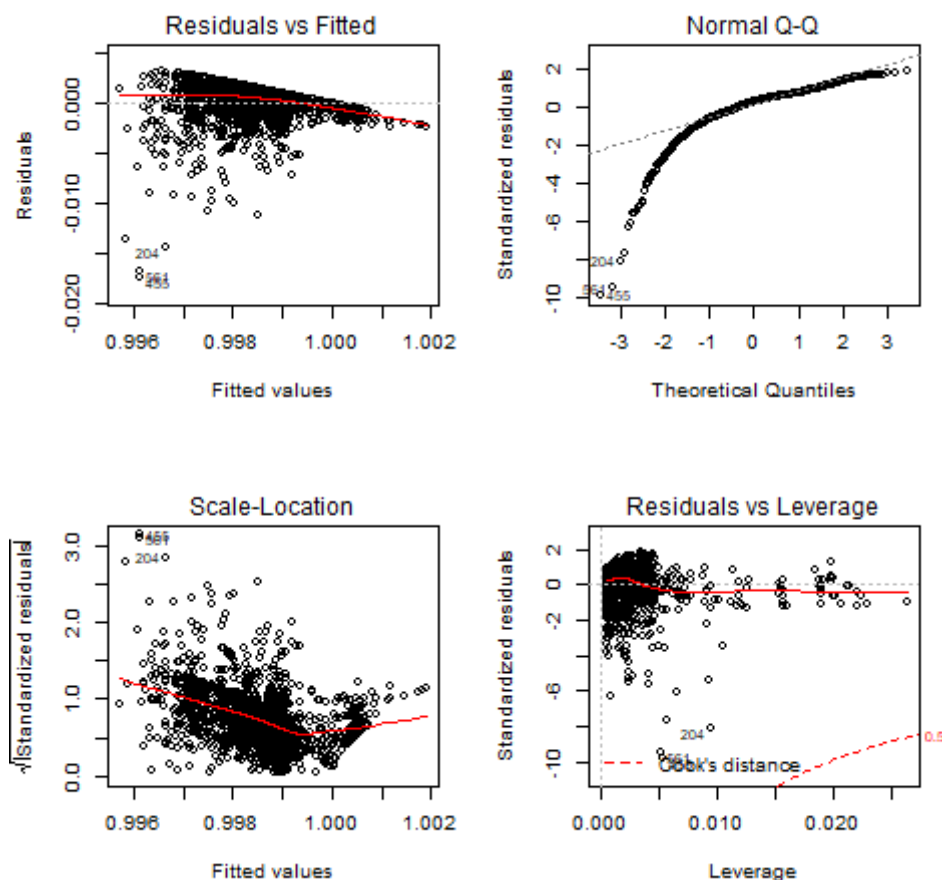


Fig 6.1 Sub Meters VS Power Factor

Except when the sub meter power is smaller than 5, the power sub meters consumed is basically linear with power factor. So, we do the multivariable regression.



```
Call:
lm(formula = build_model_2_bt_sample$clean_data.Phase ~ build_model_2_bt_sample$clean_data.Sub_metering_1 +
    build_model_2_bt_sample$clean_data.Sub_metering_2 + build_model_2_bt_sample$clean_data.Sub_metering_3)

Residuals:
    Min       1Q   Median       3Q      Max
-0.0175184 -0.0006204  0.0004065  0.0010162  0.0033315

Coefficients:
                Estimate Std. Error t value Pr(>|t|)
(Intercept)      9.938e-01  4.422e-04  2247.578 < 2e-16 ***
build_model_2_bt_sample$clean_data.Sub_metering_1  6.735e-05  4.124e-06   16.333 < 2e-16 ***
build_model_2_bt_sample$clean_data.Sub_metering_2  4.318e-05  2.831e-06   15.251 < 2e-16 ***
build_model_2_bt_sample$clean_data.Sub_metering_3  6.543e-05  2.246e-05    2.913  0.00362 **
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.001782 on 1996 degrees of freedom
Multiple R-squared:  0.188,    Adjusted R-squared:  0.1868
F-statistic: 154 on 3 and 1996 DF, p-value: < 2.2e-16
```

Fig 6.2 Fitting Multivariable Regression

We can see that the fitting is basically ok. Meter 1 contributes the most, and then the Meter 3. Meter 2 contributes the least for power factor.

### 3. Association Rules

Meter 1 represents kitchen with Dishwasher, oven and microwave. Meter with washing machine, tumble-drier, and refrigerator, Meter 3 measures the power consumption from Air Conditioner and Heater. We want to find patterns between different usage. For example, whether the



household prefer to use the washing machine after they use microwave to cook in afternoon. When we know about this thing, the owner of the household can combine the rules with the electric price. Peak and valley electric charges are applied somewhere. So, owner can find a way to adjust their custom to use the washing machine in mid night, which can lead to reduced cost.

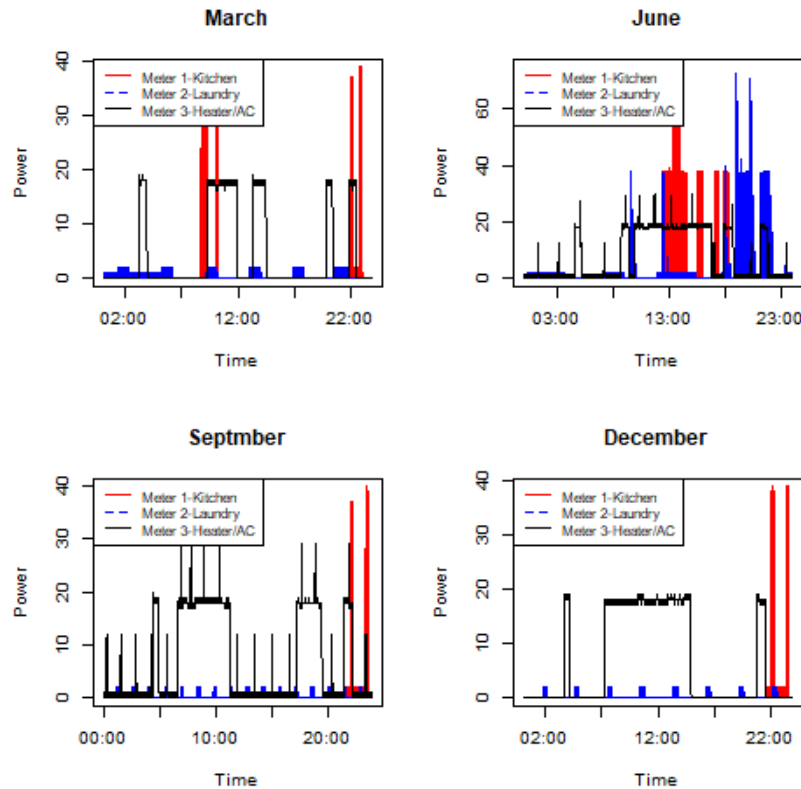


Fig 7.1 Examples of Power Consumption for a Day

We can see that different usage behaviours have different threshold value which can be reflected from the Fig 7.1. After doing some research, we can transform the data into different usage pattern with the conditions below.

Table 1.1: Conditions to transform data into proper format for Association Rule Mining

Mark	Behaviour	Power Range(W*h)	Least Duration(min)
A	Dishwasher/ Microwave	31-50	25
B	Oven	>50	45
C	Refrigerator	0-10	480
D	Washing Machine/ Drier	11-50	35
E	Washing Machine + Drier	>50	35
F	Heater	11-35	80
G	AC	16-25	480
H	Heater + AC	>25	480

Within the least duration, the same behaviour will be treated as the same. The same behaviour happens the first time in the day will be marked as 'A1', the second will be marked as 'A2'. To transform the data into the format we want. First, we find the behaviour which meets the requirements for A-H separately and store with timestamp. Then we sort all the behaviour by time,

and extract the behaviour for each day. Then we use Apriori with the conditions we set to find the rules.

```
[1] lhs rhs support confidence lift
    {A1,A2} => {F1} 0.6712707 0.9346154 1.0099426
```

Fig 7.2 One outcome from association rules mining

The Fig 7.2 means, the household usually open their Heater after they use the kitchen.

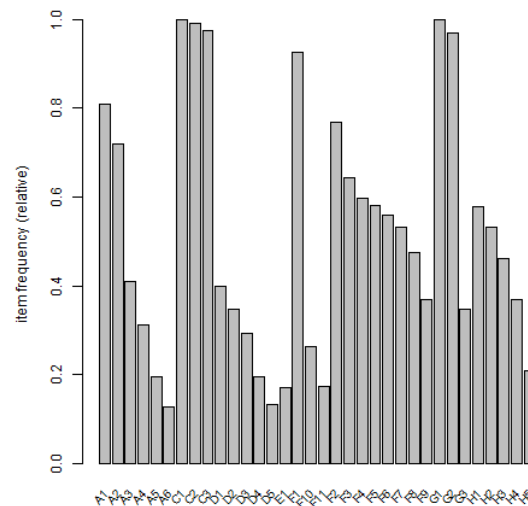


Fig 7.3 Frequency For different behaviour

The Fig 7.3 show the frequency for different behaviour, which we can find that the household didn't use washing machine frequently, but kept the refrigerator opening.

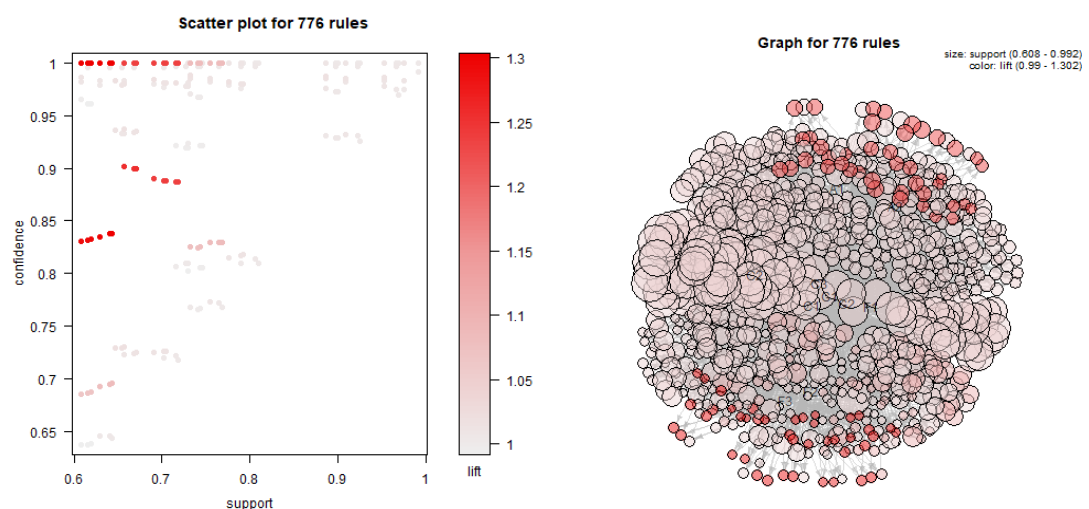


Fig 7.4 Plots for rules

## Appendix

### Code

Github Repo link: [https://github.com/PascalSun/Data\\_Science\\_Report](https://github.com/PascalSun/Data_Science_Report)

```
#####
# Prepare the working environment #
#####

Sys.setlocale(category = "LC_ALL", locale = "English")
working_dir = 'D:/Google/MIT/DataScience/Project/Power'
setwd(working_dir)
#install.packages(c("VIM","mice"))
getwd()
rm(list=ls())

#####
# Load the data from the txt file #
#####

# load data from the txt file function
load_data <- function(filename){
  Data <- read.csv(file=filename,header=TRUE,sep=";")
  return(Data)
}

# load data in mydata data.frame
raw_data <- load_data("../code/household_power_consumption.txt")
# print('Original Summary Data')
# print(summary(mydata))
raw_data$version <- 0
mydata <- raw_data

#####
# Assign Proper type to the values #
#####

# Deal with data to make them the proper data format
format_data <-function(mydata){
  mydata$DateFull <- paste(mydata$Date,mydata$Time)
  mydata$DateFull <-
as.POSIXct(mydata$DateFull,"%d/%m/%Y %H:%M:%S",tz=Sys.timezone())
  mydata$Date <- as.Date(mydata$Date,"%d/%m/%Y")
  mydata$DateTimeS <- as.POSIXct(mydata$Time,"%H:%M:%S",tz=Sys.timezone())
  mydata$Global_active_power <-
as.numeric(as.character(mydata$Global_active_power))
  mydata$Global_reactive_power <-
as.numeric(as.character(mydata$Global_reactive_power))
  mydata$Global_intensity <- as.numeric(as.character(mydata$Global_intensity))
  mydata$Voltage <- as.numeric(as.character(mydata$Voltage))
  mydata$Sub_metering_1 <- as.numeric(as.character(mydata$Sub_metering_1))
  mydata$Sub_metering_2 <- as.numeric(as.character(mydata$Sub_metering_2))
  mydata$Sub_metering_3 <- as.numeric(as.character(mydata$Sub_metering_3))
}
```

```

    return(mydata)
}

mydata = format_data(mydata)
mydata$version <- 1

#=====
=====#

#####
# Deal with missing Data #
#####

#####
# I. Analyze Missing Data Distribution #
#####

#####
# a. Summary Data #
#####
print('Summary Data')
print(summary(mydata))

#####
# b. Find Missing Rate and Counts #
#####

# Find the missing rate of the data
pMiss <- function(x) {sum(is.na(x))/length(x)*100}
print('Missing data rate for each column')
print(apply(mydata,2,pMiss))

# Analysis with mice
library(mice)
print(md.pattern(mydata))

#####
# c. Visualize the missing data #
#####

library('VIM')

# show missing data matrix
png('./img/Data-Missing-Matrix.png')
matrixplot(mydata)
dev.off()

```

```

# Data-Missing-Meter-3 check the missing values graph in sub-metering
png('./img/Data-Missing-Time-Meter-3.png')
par(mfrow=c(3,2))
marginplot(mydata[c("DateFull", "Sub_metering_3")], col=c("darkgray", "red", "blue"))
#dev.off()

#png('./img/Data-Missing-Time-Voltage.png')
marginplot(mydata[c("DateFull", "Voltage")], col=c("darkgray", "red", "blue"))
#dev.off()

#png('./img/Data-Missing-Time-Meter-2.png')
marginplot(mydata[c("DateFull", "Sub_metering_2")], col=c("darkgray", "red", "blue"))
#dev.off()

#png('./img/Data-Missing-Time-Intensity.png')
marginplot(mydata[c("DateFull", "Global_intensity")], col=c("darkgray", "red", "blue"))
#dev.off()

#png('./img/Data-Missing-Time-Global_Active.png')
marginplot(mydata[c("DateFull", "Global_active_power")], col=c("darkgray", "red", "blue"))
#dev.off()

#png('./img/Data-Missing-Time-Global_ReActive.png')
marginplot(mydata[c("DateFull", "Global_reactive_power")], col=c("darkgray", "red", "blue"))
dev.off()

#png('./img/Data-Missing-Time-Meter-3.png')
marginplot(mydata[c("DateFull", "Sub_metering_3")], col=c("darkgray", "red", "blue"))
#dev.off()

dev.off()

#####
# d. Analysis the distribution of Missing data #
#####

# The missing data can be analyzed from 2 aspects
# 1. The time is random or meaning something
# 2. The data before and after it, cause it? (High Intensity or ?)

missing_data_time = subset(mydata, is.na(mydata$DateFull))
missing_data_other = subset(mydata, is.na(mydata$Global_active_power))

```

```

#####
# Notes #
#####
# 1. The missing of the data is because some unknown of the equipment or
network, and most happen in 2007 and 2008, mainly in 3 parts.
# 2. After analyze the data near the missing data, we found that the high
voltage is not the reason why the data is missing.
# 3. we infer that it may caused because the outage nearby
# 4. So we decide to drop the data null

#####
# e. Drop the missing rows #
#####

clean_data <- na.omit(mydata)
clean_data$version <- 2

#####
# f. add other data to the column #
#####

clean_data$Sub_metering_other <- (clean_data$Global_active_power*1000/60)-
clean_data$Sub_metering_1-clean_data$Sub_metering_2-clean_data$Sub_metering_3
print('Summary for sub meter other')
print(summary(clean_data$Sub_metering_other))

# replace the sub meter other < 0 with 0
clean_data$Sub_metering_other[clean_data$Sub_metering_other < 0 ] <- 0

clean_data$version <- 3
#=====#

#####
# Explore the data #
#####

# a. summary the data
print('Summary for clean data')
print(summary(clean_data))

png('./img/Data-Range-Global_Activate_Power.png')
par(mfrow=c(4,2))
plot(clean_data$DateFull,clean_data$Global_active_power,main='Global Active
Power')
# dev.off()

# png('./img/Data-Range-Global_ReActivate_Power.png')

```

```

plot(clean_data$DateFull,clean_data$Global_reactive_power,main='Global
ReActive Power')
# dev.off()

#png('./img/Data-Range-Voltage.png')
plot(clean_data$DateFull,clean_data$Voltage,main='Voltage')
#dev.off()

#png('./img/Data-Range-Global_Intensity.png')
plot(clean_data$DateFull,clean_data$Global_intensity,main='Global Intensity')
#dev.off()

#png('./img/Data-Range-Sub_Metering_1.png')
plot(clean_data$DateFull,clean_data$Sub_metering_1,main='Sub Metering 1')
#dev.off()

#png('./img/Data-Range-Sub_Metering_2.png')
plot(clean_data$DateFull,clean_data$Sub_metering_2,main='Sub Metering 2')
#dev.off()

#png('./img/Data-Range-Sub_Metering_3.png')
plot(clean_data$DateFull,clean_data$Sub_metering_3,main='Sub Metering 3')
#dev.off()

#png('./img/Data-Range-Sub_Metering_other.png')
plot(clean_data$DateFull,clean_data$Sub_metering_other,main='Sub Metering
Other')
dev.off()

# b. Data Preprocessing
# seperate the year,month,day,hour,minute into different column
clean_data$Year = format(clean_data$DateFull,format="%Y")
clean_data$Month = format(clean_data$DateFull,format="%m")
clean_data$Day = format(clean_data$DateFull,format="%d")
clean_data$Hour = format(clean_data$DateFull,format="%H")
clean_data$Minute = format(clean_data$DateFull,format="%M")
clean_data$DayOfYear = format(clean_data$DateFull,format="%j")

# Total Power
clean_data$Global_Power <-
sqrt(clean_data$Global_active_power^2+clean_data$Global_reactive_power^2)
clean_data$Phase <- clean_data$Global_active_power/clean_data$Global_Power

png('./img/Data-Explore-Phase-VS-Voltage.png')
par(mfrow=c(3,1))
plot(clean_data$Phase,clean_data$Voltage)
#dev.off()

```

```

#png('./img/Data-Explore-Phase-VS-Intensity.png')
plot(clean_data$Phase,clean_data$Global_intensity)
#dev.off()

# png('./img/Data-Explore-Global_Power-vs-Date.png')
# plot(clean_data$DateFull,clean_data$Global_Power)
# dev.off()

#png('./img/Data-Explore-Voltage-vs-Intensity.png')
plot(clean_data$Voltage,clean_data$Global_intensity)
dev.off()

png('./img/Data-Explore-Pie-Chart-Meters.png')
slice <-
c(sum(clean_data$Sub_metering_1),sum(clean_data$Sub_metering_2),sum(clean_data
$Sub_metering_3),sum(clean_data$Sub_metering_other))
lbls <- c('Sub_Meter_1','Sub_Meter_2','Sub_Meter_3','Sub_Meter_Other')
piepercent<- round(100*slice/sum(slice), 1)
pie(slice,labels = piepercent,main="Pie Chart of the Power
Consumption",col=rainbow(length(slice)))
legend("topright", c("Meter1","Meter2","Meter3","Meter_other"), cex = 0.8,fill
= rainbow(length(slice)))
dev.off()

#####
# Build Model 1: Power Consumption VS Month #
#####

#####
# For Days #
#####

# 2007
year2007 = subset(clean_data,clean_data$Year=='2007')
year2007_Power <-
aggregate(year2007$Global_Power,by=list(DayofYear=year2007$DayOfYear),FUN=sum)

png('./img/Data-Build-Model-regression-1-check-2007.png')
par(mfrow=c(3,2))
plot(year2007_Power$DayofYear,year2007_Power$x,xlab="DayofYear",ylab="Power",m
ain='Consumption of Power the whole year 2007')
#dev.off()

# 2008

```



```

year2008 = subset(clean_data,clean_data$Year=='2008')
year2008_Power <-
aggregate(year2008$Global_Power,by=list(DayofYear=year2008$DayOfYear),FUN=sum)

#png('./img/Data-Build-Model-regression-1-check-2008.png')
plot(year2008_Power$DayofYear,year2008_Power$x,xlab="DayofYear",ylab="Power",m
ain='Consumption of Power the whole year 2008')
#dev.off()

# 2009
year2009 = subset(clean_data,clean_data$Year=='2009')
year2009_Power <-
aggregate(year2009$Global_Power,by=list(DayofYear=year2009$DayOfYear),FUN=sum)

#png('./img/Data-Build-Model-regression-1-check-2009.png')
plot(year2009_Power$DayofYear,year2009_Power$x,xlab="DayofYear",ylab="Power",m
ain='Consumption of Power the whole year 2009')
#dev.off()

#####
# For month #
#####

# 2007
year2007 = subset(clean_data,clean_data$Year=='2007')
year2007_Power_Month <-
aggregate(year2007$Global_Power,by=list(Month=year2007$Month),FUN=sum)

#png('./img/Data-Build-Model-regression-1-check-2007-month.png')
#par(mfrow=c(3,1))
plot(year2007_Power_Month$Month,year2007_Power_Month$x,xlab="Month",ylab="Powe
r",main='Consumption of Power the whole year 2007')
#dev.off()

# 2008
year2008 = subset(clean_data,clean_data$Year=='2008')
year2008_Power_Month <-
aggregate(year2008$Global_Power,by=list(Month=year2008$Month),FUN=sum)

#png('./img/Data-Build-Model-regression-1-check-2008-month.png')
plot(year2008_Power_Month$Month,year2008_Power_Month$x,xlab="Month",ylab="Powe
r",main='Consumption of Power the whole year 2008')
#dev.off()

# 2009
year2009 = subset(clean_data,clean_data$Year=='2009')
year2009_Power_Month <-
aggregate(year2009$Global_Power,by=list(Month=year2009$Month),FUN=sum)

```

```

#png('./img/Data-Build-Model-regression-1-check-2009-month.png')
plot(year2009_Power_Month$Month,year2009_Power_Month$x,xlab="Month",ylab="Power",main='Consumption of Power the whole year 2009')
dev.off()

#####
# Box Plot #
#####

png('./img/Data-Build-Model-regression-1-boxplot-2007-day.png')
par(mfrow=c(3,2))
boxplot(year2007_Power$x,main="Year2007Power")
# dev.off()

# png('./img/Data-Build-Model-regression-1-boxplot-2007-Month.png')
boxplot(year2007_Power_Month$x,main="Year2007PowerMonth")
# dev.off()

# png('./img/Data-Build-Model-regression-1-boxplot-2008-day.png')
boxplot(year2008_Power$x,main="Year2008Power")
# dev.off()

# png('./img/Data-Build-Model-regression-1-boxplot-2008-Month.png')
boxplot(year2008_Power_Month$x,main="Year2008PowerMonth")
# dev.off()

# png('./img/Data-Build-Model-regression-1-boxplot-2009-day.png')
boxplot(year2009_Power$x,main="Year2009Power")
# dev.off()

# png('./img/Data-Build-Model-regression-1-boxplot-2009-Month.png')
boxplot(year2009_Power_Month$x,main="Year2009PowerMonth")
dev.off()

#####
# Model Build #
#####

year2007_Power_Month$Month <-
as.numeric(as.character(year2007_Power_Month$Month))
m <-
nls(x~a*Month^2+b*Month+c,data=year2007_Power_Month,start=list(a=300,b=11,c=7)
)

```

```

png('./img/Data-Build-Model-regression-1-Model-and-Validate.png')
par(mfrow=c(3,1))
plot(year2007_Power_Month$Month,year2007_Power_Month$x,xlab="Month(2007)",ylab="Power")
lines(year2007_Power_Month$Month,predict(m),col="red",lty=2,lwd=3)

legend('top',legend=capture.output(summary(m)),cex=0.6)
#dev.off()

#png('./img/Data-Build-Model-regression-1-Model-Validate-2008.png')
plot(year2008_Power_Month$Month,year2008_Power_Month$x,xlab="Month(2008)",ylab="Power")
lines(year2008_Power_Month$Month,predict(m),col="red",lty=2,lwd=3)
#dev.off()

#png('./img/Data-Build-Model-regression-1-Model-Validate-2009.png')
plot(year2009_Power_Month$Month,year2009_Power_Month$x,xlab="Month(2009)",ylab="Power")
lines(year2009_Power_Month$Month,predict(m),col="red",lty=2,lwd=3)
dev.off()

#####
# Build Model 2: Regression between Rate and sub meters #
#####

#####
# Observation #
#####

build_model_2 <-
data.frame(clean_data$Phase,clean_data$Sub_metering_1,clean_data$Sub_metering_
2,clean_data$Sub_metering_3)

png('./img/Data-Build-Model-regression-2-Sub-Meter-1-VS-Rate.png')
par(mfrow=c(3,1))
plot(build_model_2$clean_data.Sub_metering_1,build_model_2$clean_data.Phase,xl
ab="Sub Meter 1",ylab="Sub Meter 2")
# dev.off()

# png('./img/Data-Build-Model-regression-2-Sub-Meter-2-VS-Rate.png')
plot(build_model_2$clean_data.Sub_metering_2,build_model_2$clean_data.Phase,xl
ab="Sub Meter 1",ylab="Sub Meter 2")
# dev.off()

# png('./img/Data-Build-Model-regression-2-Sub-Meter-3-VS-Rate.png')
plot(build_model_2$clean_data.Sub_metering_3,build_model_2$clean_data.Phase,xl
ab="Sub Meter 1",ylab="Sub Meter 2")
dev.off()

```

```

# we found that when the value is larger than 5, it seems like a linear, so we
drop the small values

#####
# Sampling #
#####

build_model_2_bt =
subset(build_model_2,(build_model_2$clean_data.Sub_metering_1>5)&(build_model_
2$clean_data.Sub_metering_2>5)&(build_model_2$clean_data.Sub_metering_3>5))

build_model_2_bt_sample <-
build_model_2_bt[sample(nrow(build_model_2_bt),2000),]

build_model_2_test <- build_model_2_bt[sample(nrow(build_model_2_bt),2000),]

#####
# Build Model #
#####

fit <-
lm(build_model_2_bt_sample$clean_data.Phase~build_model_2_bt_sample$clean_data
.Sub_metering_1+build_model_2_bt_sample$clean_data.Sub_metering_2+build_model_
2_bt_sample$clean_data.Sub_metering_3)
print(summary(fit))
png('./img/Data-Build-Model-regression-2-Report.png')
par(mfrow=c(2,2))
plot(fit)
dev.off()

png('./img/Data-Build-Model-regression-2-Fitted.png')
plot(build_model_2_bt_sample$clean_data.Phase,fitted(fit))
dev.off()

pred <- predict(fit,build_model_2_test)
png('./img/Data-Build-Model-regression-2-Test.png')
plot(pred,build_model_2_test$clean_data.Phase-pred)
dev.off()

#####
# Build Model 3: Association Rules between different meters #
#####
# Packages: arules arulesViz

# 1. Sub Meter 1: Kitchen: Dishwasher, Oven, Microwave

```

```

# 2. Sub Meter 2: Laundry: Washing Machine, Tumble-Drier, Refrigerator, a
light
# 3. Sub Meter 3: AC/Heater

#####
# Example Plot for a day #
#####

png('./img/Data-Build-Model-Rule-1-Day-Example-Sub-Meter.png')
par(mfrow=c(2,2))
day2008_01_01 <- subset(clean_data,clean_data$Date == as.Date('2008-03-15'))
plot(day2008_01_01$DateFull,day2008_01_01$Sub_metering_1,col="red",xlab="Time"
,ylab="Power",type="h",main="March")
lines(day2008_01_01$DateFull,day2008_01_01$Sub_metering_2,type="h",col="blue")
lines(day2008_01_01$DateFull,day2008_01_01$Sub_metering_3,type="l",col="black"
)
legend('topleft',legend = c('Meter 1-Kitchen','Meter 2-Laundry', 'Meter 3-
Heater/AC'),col=c('red','blue','black'),cex=0.8,lty=1:2)

day2008_01_01 <- subset(clean_data,clean_data$Date == as.Date('2008-06-15'))
plot(day2008_01_01$DateFull,day2008_01_01$Sub_metering_1,col="red",xlab="Time"
,ylab="Power",type="h",main="June")
lines(day2008_01_01$DateFull,day2008_01_01$Sub_metering_2,type="h",col="blue")
lines(day2008_01_01$DateFull,day2008_01_01$Sub_metering_3,type="l",col="black"
)
legend('topleft',legend = c('Meter 1-Kitchen','Meter 2-Laundry', 'Meter 3-
Heater/AC'),col=c('red','blue','black'),cex=0.8,lty=1:2)

day2008_01_01 <- subset(clean_data,clean_data$Date == as.Date('2008-09-15'))
plot(day2008_01_01$DateFull,day2008_01_01$Sub_metering_1,col="red",xlab="Time"
,ylab="Power",type="h",main="Septmber")
lines(day2008_01_01$DateFull,day2008_01_01$Sub_metering_2,type="h",col="blue")
lines(day2008_01_01$DateFull,day2008_01_01$Sub_metering_3,type="l",col="black"
)
legend('topleft',legend = c('Meter 1-Kitchen','Meter 2-Laundry', 'Meter 3-
Heater/AC'),col=c('red','blue','black'),cex=0.8,lty=1:2)

day2008_01_01 <- subset(clean_data,clean_data$Date == as.Date('2008-12-15'))
plot(day2008_01_01$DateFull,day2008_01_01$Sub_metering_1,col="red",xlab="Time"
,ylab="Power",type="h",main="December")
lines(day2008_01_01$DateFull,day2008_01_01$Sub_metering_2,type="h",col="blue")
lines(day2008_01_01$DateFull,day2008_01_01$Sub_metering_3,type="l",col="black"
)
legend('topleft',legend = c('Meter 1-Kitchen','Meter 2-Laundry', 'Meter 3-
Heater/AC'),col=c('red','blue','black'),cex=0.8,lty=1:2)

dev.off()

```

```

png('./img/Data-Build-Model-Rule-1-Day-Example-Power.png')
plot(day2008_01_01$DateFull,day2008_01_01$Global_active_power,xlab="Time",ylab=
="Active Power",type='h',main="Activate Power of A Day")
dev.off()

#####
# Translate into Patterns for days #
#####

# a. Store Time and Way into pattern
pattern <- data.frame('Time'=character(),'Mark'=character(),stringsAsFactors =
FALSE)

Month_March <- subset(clean_data,(clean_data$Date > as.Date('2008-11-
07')&(clean_data$Date< as.Date('2009-11-06'))))

# A: When use Dishwasher or Microwave in kitchen, sub_metering_1:31~50,
classfiy A

A <-
subset(Month_March,Month_March$Sub_metering_1>30&Month_March$Sub_metering_1<=5
0)

start_time = A[1,]$DateFull
end_time = A[nrow(A),]$DateFull
n <- 2

while(n<nrow(A)){
  n <- n+1
  if(difftime(A[n,]$DateFull,start_time,units="mins")<25){
    n <- n+1
  }else{
    # add it to column here
    print(start_time)
    pattern[nrow(pattern)+1,] <- c(as.character(start_time),'A')
    start_time <- A[n,]$DateFull
  }
}

# B: When use oven in kitchen, sub_metering_1: >50, classify B
B <- subset(Month_March,Month_March$Sub_metering_1>50)

start_time = B[1,]$DateFull
end_time = B[nrow(B),]$DateFull
n <- 2

while(n<nrow(B)){

```

```

n <- n+1
if(difftime(B[n,]$DateFull,start_time,units="mins")<45){
  n <- n+1
}else{
  # add it to column here
  print(start_time)
  pattern[nrow(pattern)+1,] <- c(as.character(start_time),'B')
  start_time <- B[n,]$DateFull
}
}

# C: When refrigerator or light is used in laundry, sub_metering_2 <= 10

C <- subset(Month_March,Month_March$Sub_metering_2<=10)

start_time = C[1,]$DateFull
print(start_time)
end_time = C[nrow(C),]$DateFull
n <- 2

while(n<nrow(C)){
  n <- n+1
  if(difftime(C[n,]$DateFull,start_time,units="mins")<480){
    n <- n+1
  }else{
    # add it to column here
    print(start_time)
    pattern[nrow(pattern)+1,] <- c(as.character(start_time),'C')
    start_time <- C[n,]$DateFull
  }
}

# D: When refrigerator or light and washing-machine or a tumble-drier is used
in laundry, sub_metering_2: 10~50

D <-
subset(Month_March,(Month_March$Sub_metering_2>10)&(Month_March$Sub_metering_2
<=50))

start_time = D[1,]$DateFull
end_time = D[nrow(D),]$DateFull
n <- 2

while(n<nrow(D)){
  n <- n+1
  if(difftime(D[n,]$DateFull,start_time,units="mins")<35){

```

```

    n <- n+1
  }else{
    # add it to column here
    print(start_time)
    pattern[nrow(pattern)+1,] <- c(as.character(start_time),'D')
    start_time <- D[n,]$DateFull
  }
}

# E: When refrigerator AND washing-machine, a tumble-drier is used in
laundry, sub_metering_2: 50

E <- subset(Month_March,Month_March$Sub_metering_2>50)

start_time = E[1,]$DateFull
end_time = E[nrow(E),]$DateFull
n <- 2

while(n<nrow(E)){
  n <- n+1
  if(difftime(E[n,]$DateFull,start_time,units="mins")<35){
    n <- n+1
  }else{
    # add it to column here
    print(start_time)
    pattern[nrow(pattern)+1,] <- c(as.character(start_time),'E')
    start_time <- E[n,]$DateFull
  }
}

# F: heater working 11-15

F <-
subset(Month_March,(Month_March$Sub_metering_3)>10&(Month_March$Sub_metering_3
)<=15)

start_time = F[1,]$DateFull
end_time = F[nrow(F),]$DateFull
n <- 2

while(n<nrow(F)){
  n <- n+1
  if(difftime(F[n,]$DateFull,start_time,units="mins")<80){
    n <- n+1
  }else{
    # add it to column here

```



```

    print(start_time)
    pattern[nrow(pattern)+1,] <- c(as.character(start_time),'F')
    start_time <- F[n,]$DateFull
  }
}

# G: AC 16-25

G <-
subset(Month_March,(Month_March$Sub_metering_3)>15&(Month_March$Sub_metering_3
)<=25)

start_time = G[1,]$DateFull
end_time = G[nrow(G),]$DateFull
n <- 2

while(n<nrow(G)){
  n <- n+1
  if(difftime(G[n,]$DateFull,start_time,units="mins")<480){
    n <- n+1
  }else{
    # add it to column here
    print(start_time)
    pattern[nrow(pattern)+1,] <- c(as.character(start_time),'G')
    start_time <- G[n,]$DateFull
  }
}

# H: AC+Heater >25

H <- subset(Month_March,Month_March$Sub_metering_3>25)

start_time = H[1,]$DateFull
end_time = H[nrow(H),]$DateFull
n <- 2

while(n<nrow(H)){
  n <- n+1
  if(difftime(H[n,]$DateFull,start_time,units="mins")<80){
    n <- n+1
  }else{
    # add it to column here
    print(start_time)
    pattern[nrow(pattern)+1,] <- c(as.character(start_time),'H')
    start_time <- H[n,]$DateFull
  }
}

```

```

# b. sort pattern to become day pattern

# process to be sorted
pattern$Mark <- as.factor(pattern$Mark)
pattern$Time <-
as.POSIXct(pattern$Time,format="%Y-%m-%d %H:%M:%S",tz="Australia/Perth")
pattern$Date <- as.Date(pattern$Time,"%Y-%m-%d %H:%M:%S")
# sort by time
pattern <- pattern[order(pattern$Time),]
print('here')
# add to vector by date

n =1

this_time <- pattern[n,]$Date
this_vector <- c()
this_list <- list()
while(n < nrow(pattern)-1){
  print(this_time)
  print(n)
  if(!is.na(pattern[n,]$Date) & this_time == pattern[n,]$Date){
    # add the value into vector
    start_mark <-1

    while(paste0(pattern[n,]$Mark,as.character(start_mark)) %in% this_vector){
      start_mark <- start_mark+1
    }
    this_vector <-
c(this_vector,as.character(paste0(pattern[n,]$Mark,as.character(start_mark))))

    n <- n+1
  }else if(is.na(pattern[n,]$Date)){
    n<- n+1
  }

  else{
    key <- as.character(this_time)
    this_list[[key]] <- this_vector
    this_vector <- c()
    this_time <-pattern[n,]$Date
  }
}

#####
# Apriori to find Association Rules #

```

```
#####  
library(arules)  
trData <- as(this_list, 'transactions')  
inspect(trData)  
  
png('./img/Data-Build-Model-Rule-1-Association-Rules-Frequency.png')  
itemFrequencyPlot(trData,support=0.1,cex.names=0.8)  
dev.off()  
  
rules <- apriori(trData, parameter = list(supp = 0.6, conf  
=0.6,minlen=3,maxlen=4))  
inspect(rules)  
  
library(arulesViz)  
png('./img/Data-Build-Model-Rule-1-Rules-Scatterplot.png')  
plot(rules, method='scatterplot')  
dev.off()  
  
png('./img/Data-Build-Model-Rule-1-Rules-Graph.png')  
plot(rules, method='graph', control = list(type='items'))  
dev.off()
```