

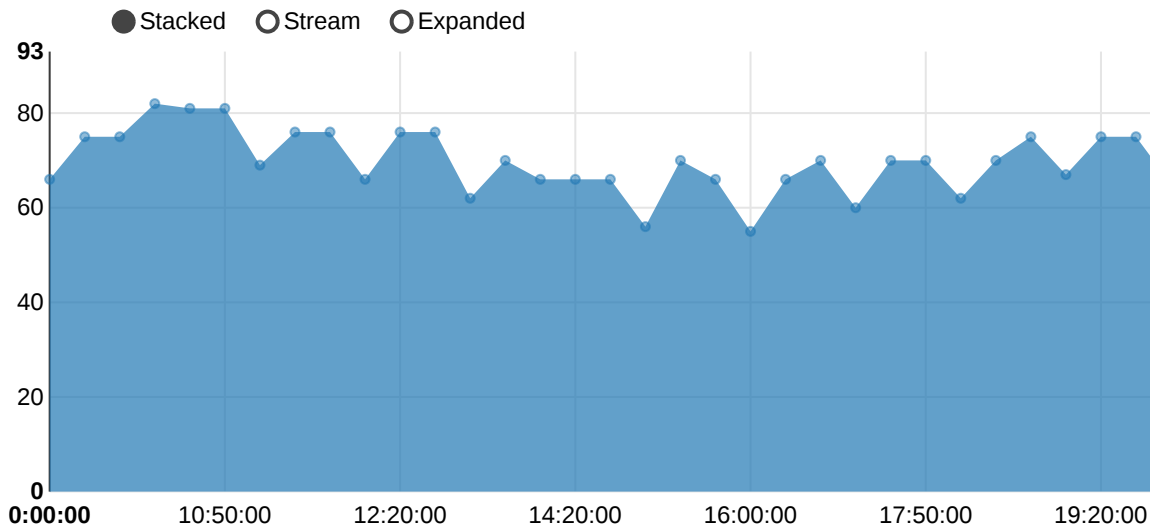
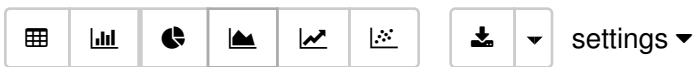
```
val humData = sqlContext.read.format("com.databricks.spark.csv").load("/home/scalaface/Desktop/sen-  
humData.toDF().registerTempTable("humidity")
```

humData: org.apache.spark.sql.Dataset[org.apache.spark.sql.Row] = [_c0: string, _c1: string ... 1 more field]

warning: there was one deprecation warning; re-run with -deprecation for details

Took 15 sec. Last updated by anonymous at April 18 2017, 7:21:58 PM.

```
%sql  
select _c0, _c1 as time,SUM(_c2) as humidity_per_day from humidity where _c0 = '13-02-2014' group by time
```



Took 5 sec. Last updated by anonymous at April 18 2017, 7:24:11 PM. (outdated)

READY    

```
import org.apache.spark.sql.functions._  
import org.joda.time.format.DateTimeFormat  
import org.apache.commons.io.IOUtils
```

FINISHED    

Zeppelin Visualiza

Navigation icons: play, expand, book, tablet, copy, download, print, trash, search.

anonymou ▾




 default ▼

FINISHED

FINISHED    

FINISHED ▶ 🔍 📖 ⚙️

↓

Took 0 sec. Last updated by anonymous at April 18 2017, 6:16:38 PM.



Zeppelin Notebook

Zeppelin Visualiza...

```
%pyspark
head(weather.ix[:,1:])
```

ERROR anonymous

```
%pyspark
weather.ix[:,1:].sum(axis=1)
```

FINISHED default

```
%pyspark
weather.ix[:,1:].mean(axis=1,skipna=False)
```

FINISHED

```
%pyspark
weather.isnull().any()
```

FINISHED

```
0                False
2014-02-13 00:00:00  False
-0                True
2014-02-13 00:00:00  False
2014-02-13 00:00:00  False
2014-02-13 00:00:00  False
2014-02-13 00:00:00  False
dtype: bool
```

Took 0 sec. Last updated by anonymous at April 18 2017, 6:14:42 PM.

```
%pyspark
#In this step, we will try to update null values.
#filling null values could be complicated.As we seen in previous data exploration steps
#that 116 was the maximum null values and total datasize is 12563. Since, maximum percent of null v
#So, null values will be replaced by mean of the particular parameter.
def updatenullvalues(dataset):
    for col in dataset.ix[:,1:]:
        if dataset[col].isnull().any:
            mean = dataset[col].mean()
            dataset[col].fillna(mean,inplace=True)
    return dataset
```

FINISHED

Took 0 sec. Last updated by anonymous at April 18 2017, 6:14:48 PM.

```
%pyspark

#Let's update null values in our dataset.
weather = updatenullvalues(weather_dataset)
#verify is there still any null value left in the dataset
weather.isnull().sum()
```

ERROR



Zeppelin Notebook



Zeppelin Visualiza...

```
import matplotlib.pyplot as plt
```

Now we are in good state as our null values are vanished.
In this code step, we will check distribution of our data.



anonymot ▾

```
data = [weather.ix[:,1], weather.ix[:,2], weather.ix[:,3], weather.ix[:,4], weather.ix[:,5]]
```

```
parameter_names = ['Dewpoint', 'Humidity', 'Pressure', 'Temperature', 'WindDirection']
```

```
fig, axis = plt.subplots()
axis.set_title("Distribution of weather parameters")
axis.set_xlabel('Weather Parameters')
axis.set_ylabel('Values')
day_plot = plt.boxplot(data, sym='o', vert=1, whis=1.5)
plt.setp(day_plot['boxes'], color = 'black')
plt.setp(day_plot['whiskers'], color = 'black')
plt.setp(day_plot['fliers'], color = 'black', marker = 'o')
axis.set_xticklabels(parameter_names)
plt.show()
```

READY ▶ ⌵ ⌵ ⌵ ⌵

```
%pyspark
# column-wise and Multiple Function Application
grouped_pressure = weather.groupby(['Pressure'])
```

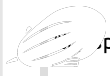
READY ▶ ⌵ ⌵ ⌵ ⌵

```
%pyspark
grouped_pressure.apply(lambda weather: weather['Humidity'].corr(weather['Temperature']))
```

READY ▶ ⌵ ⌵ ⌵ ⌵

```
%pyspark
import statsmodels.api as sm
def regression(data, yvar, xvars):
    Y = data[yvar]
    X = data[xvars]
    X['intercept'] = 1.
    result = sm.OLS(Y,X).fit()
    return result.params
```

READY ▶ ⌵ ⌵ ⌵ ⌵



Zeppelin Notebook

Zeppelin Visualiza

```
#regression(weather_dataset,'Humidity','Temperature')
```

```
grouped_by(pressure) %>% summarise(regression = lm(Humidity ~ Temperature, data = weather_dataset[pressure == pressure, 'Humidity', 'Temperature', 'Winddirection']))
```



anonymot

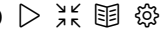


default

```
%spark.r
```

```
dewpoint = read.csv("/home/scarface/Desktop/sem-3/capstone/sravan/raw_weather_data_aarhus/dewptm.csv")
humidity = read.csv("/home/scarface/Desktop/sem-3/capstone/sravan/raw_weather_data_aarhus/hum.csv", header = TRUE)
pressure = read.csv("/home/scarface/Desktop/sem-3/capstone/sravan/raw_weather_data_aarhus/pressurem.csv", header = TRUE)
temp = read.csv("/home/scarface/Desktop/sem-3/capstone/sravan/raw_weather_data_aarhus/tempm.csv", header = TRUE)
wind = read.csv("/home/scarface/Desktop/sem-3/capstone/sravan/raw_weather_data_aarhus/wdird.csv", header = TRUE)
```

FINISHED

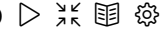


Took 0 sec. Last updated by anonymous at April 18 2017, 6:35:55 PM.

```
%spark.r
```

```
dataSet <- cbind(dewpoint, humidity$V3, pressure$V3, temp$V3, wind$V3)
```

FINISHED

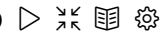


Took 1 sec. Last updated by anonymous at April 18 2017, 7:51:18 PM.

```
%spark.r
```

```
colnames(dataSet) <- c("id", "date", "dew", "hum", "pre", "tem", "win")
head(dataSet)
dataSet <- na.omit(dataSet)
```

FINISHED



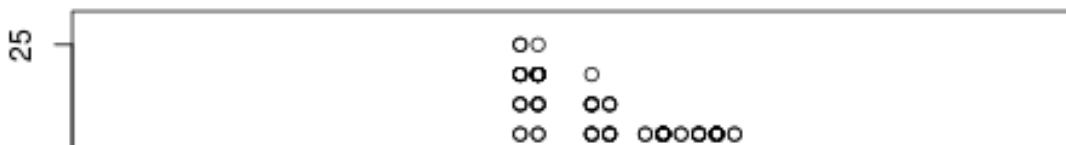
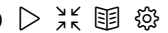
```
id      date      dew  hum  pre  tem  win
1 0 2014-02-13 00:00:00 0 66 995 4 160
2 1 2014-02-13 00:20:00 0 75 994 4 160
3 2 2014-02-13 00:50:00 0 75 993 4 160
4 3 2014-02-13 01:00:00 0 65 994 4 160
5 4 2014-02-13 01:20:00 0 75 993 4 170
6 5 2014-02-13 01:50:00 0 75 993 4 170
```

Took 0 sec. Last updated by anonymous at April 18 2017, 7:55:17 PM.

```
%spark.r
```

```
plot(dataSet[5:6])
```

FINISHED

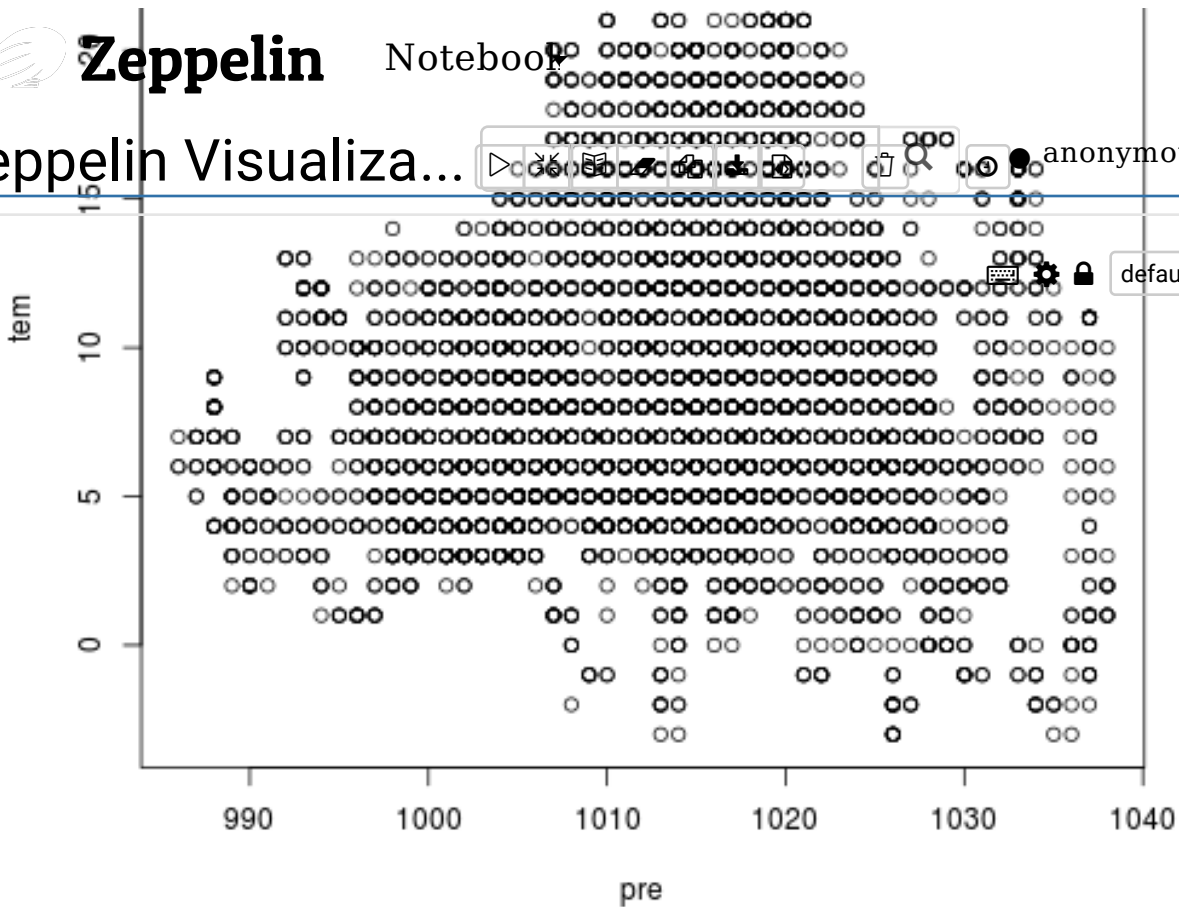




Zeppelin

Notebook

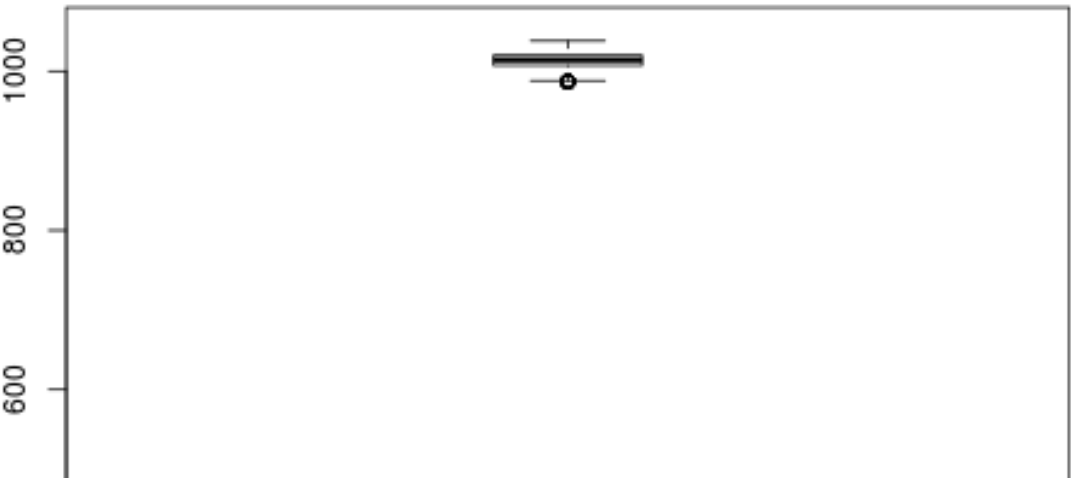
Zeppelin Visualiza...



Took 0 sec. Last updated by anonymous at April 18 2017, 7:51:25 PM.

```
%spark.r  
boxplot(dataSet[3:7])
```

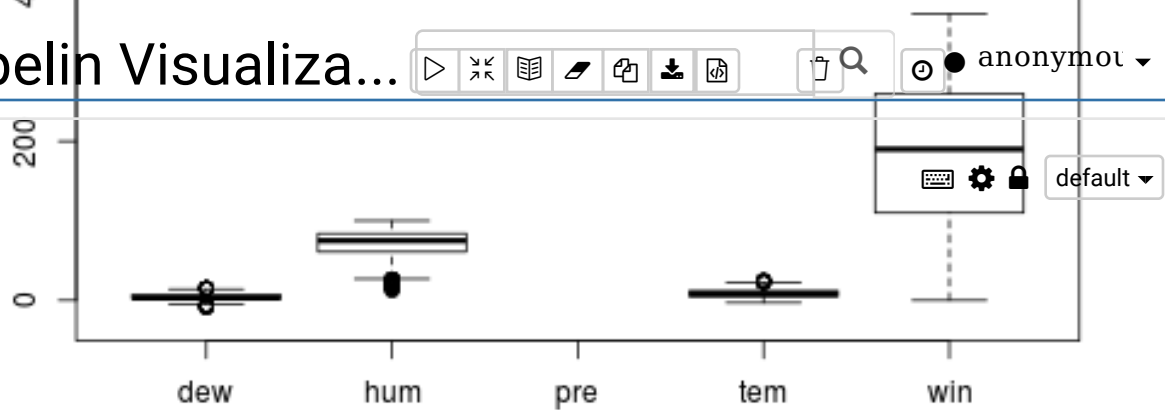
FINISHED ▶ ⌵ ⌵ ⌵ ⌵





Zeppelin Notebook

Zeppelin Visualiza...



Took 0 sec. Last updated by anonymous at April 18 2017, 7:34:45 PM.

```
%spark.r
ds=max(na.omit(dataSet[,3]))
hs=max(na.omit(dataSet[,4]))
ps=max(na.omit(dataSet[,5]))
ts=max(na.omit(dataSet[,6]))
ws=max(na.omit(dataSet[,7]))

dataSet<-na.omit(dataSet)
```

FINISHED ▶ 🔍 📖 ⚙️

Took 0 sec. Last updated by anonymous at April 18 2017, 7:36:19 PM.

```
%spark.r
for (i in 1:length(dataSet$id)){
  dataSet[i,3]<-dataSet[i,3]/ds
  dataSet[i,4]<-dataSet[i,4]/hs
  dataSet[i,5]<-dataSet[i,5]/ps
  dataSet[i,6]<-dataSet[i,6]/ts
  dataSet[i,7]<-dataSet[i,7]/ws
}
#dataSet
```

FINISHED ▶ 🔍 📖 ⚙️

Took 4 sec. Last updated by anonymous at April 18 2017, 7:39:38 PM.

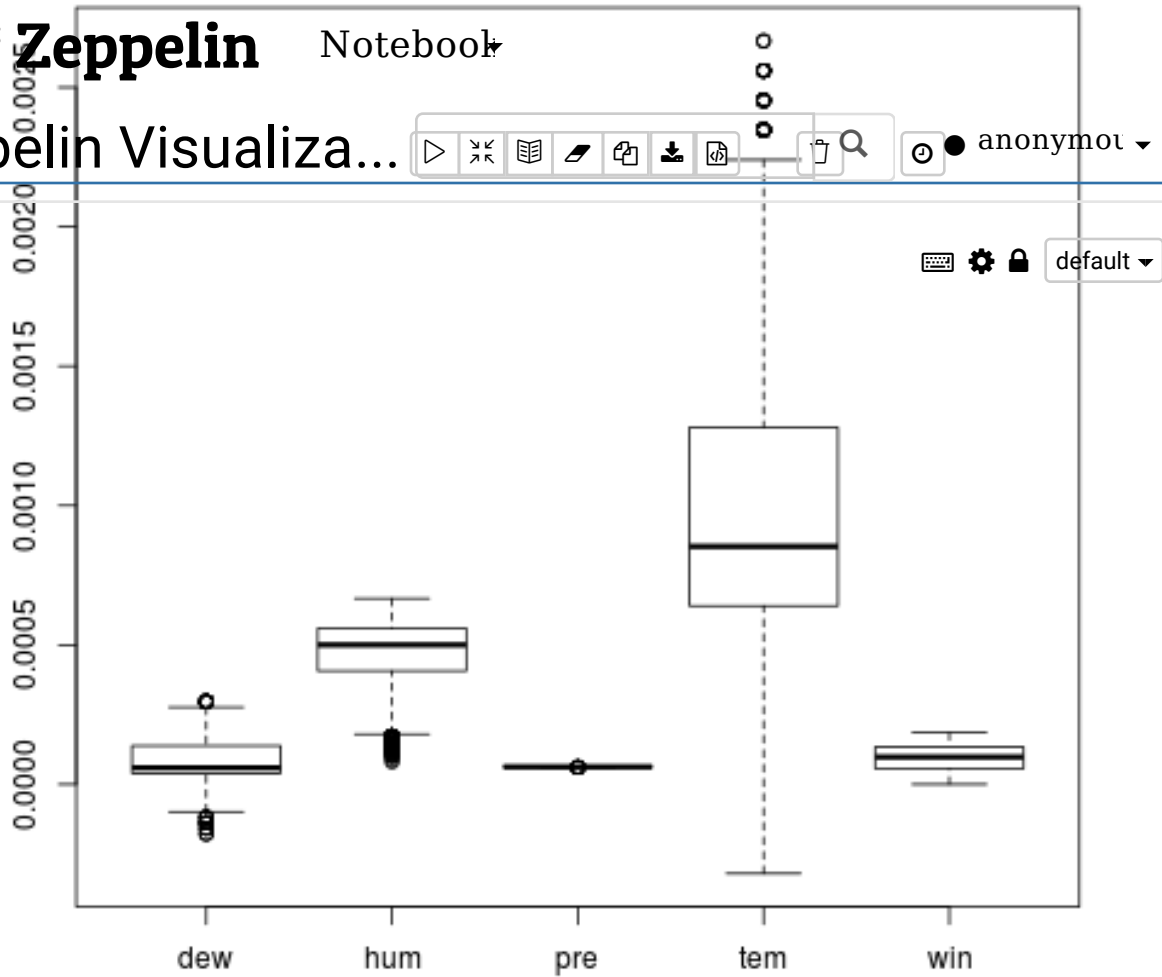
```
%spark.r
boxplot(dataSet[3:7])
```

FINISHED ▶ 🔍 📖 ⚙️



Zeppelin Notebook

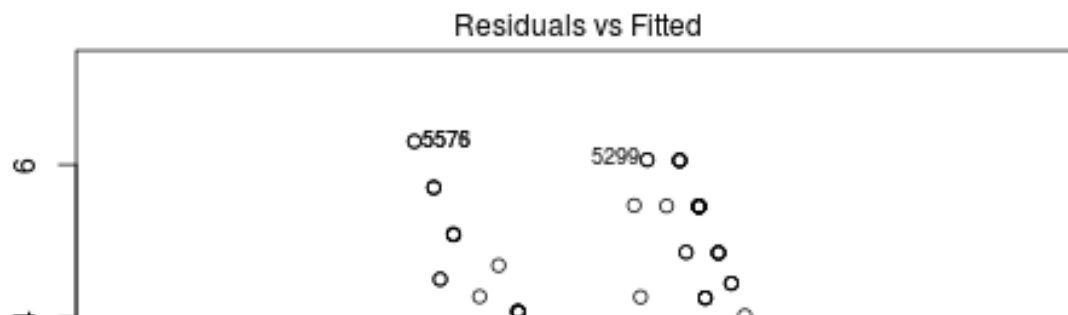
Zeppelin Visualiza...

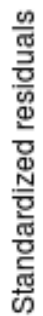
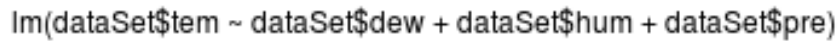


Took 0 sec. Last updated by anonymous at April 18 2017, 7:40:07 PM.

```
%spark.r
lmF <- lm(dataSet$tem ~ dataSet$dew+dataSet$hum+dataSet$pre)
plot(lmF)
```

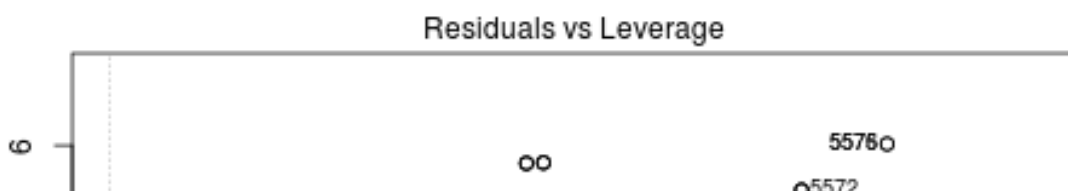
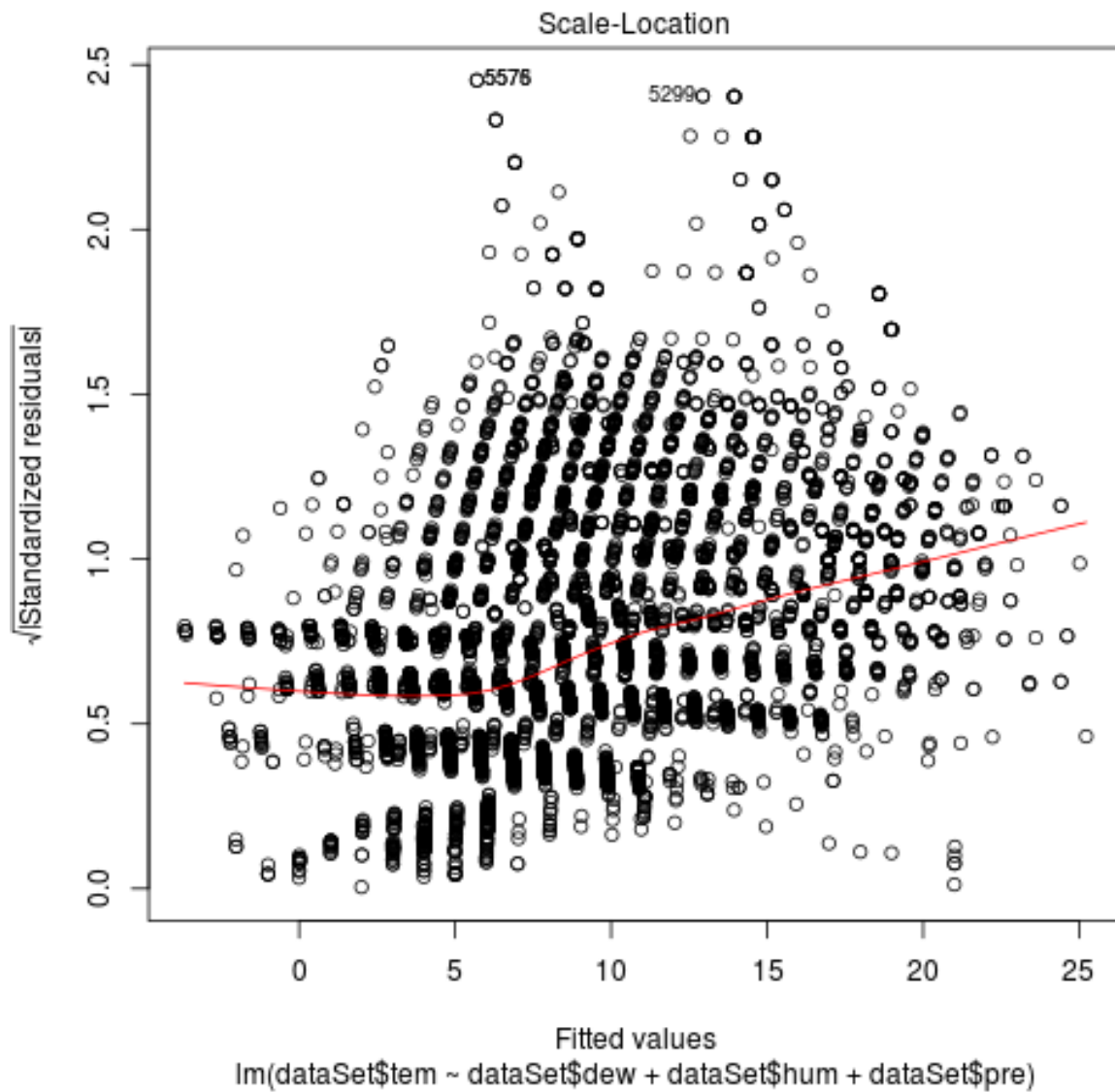
FINISHED ▶ 🔍 📖 ⚙️





Theoretical Quantiles
 $\text{lm}(\text{dataSet}\$stem \sim \text{dataSet}\$dew + \text{dataSet}\$hum + \text{dataSet}\$pre)$

   default ▾

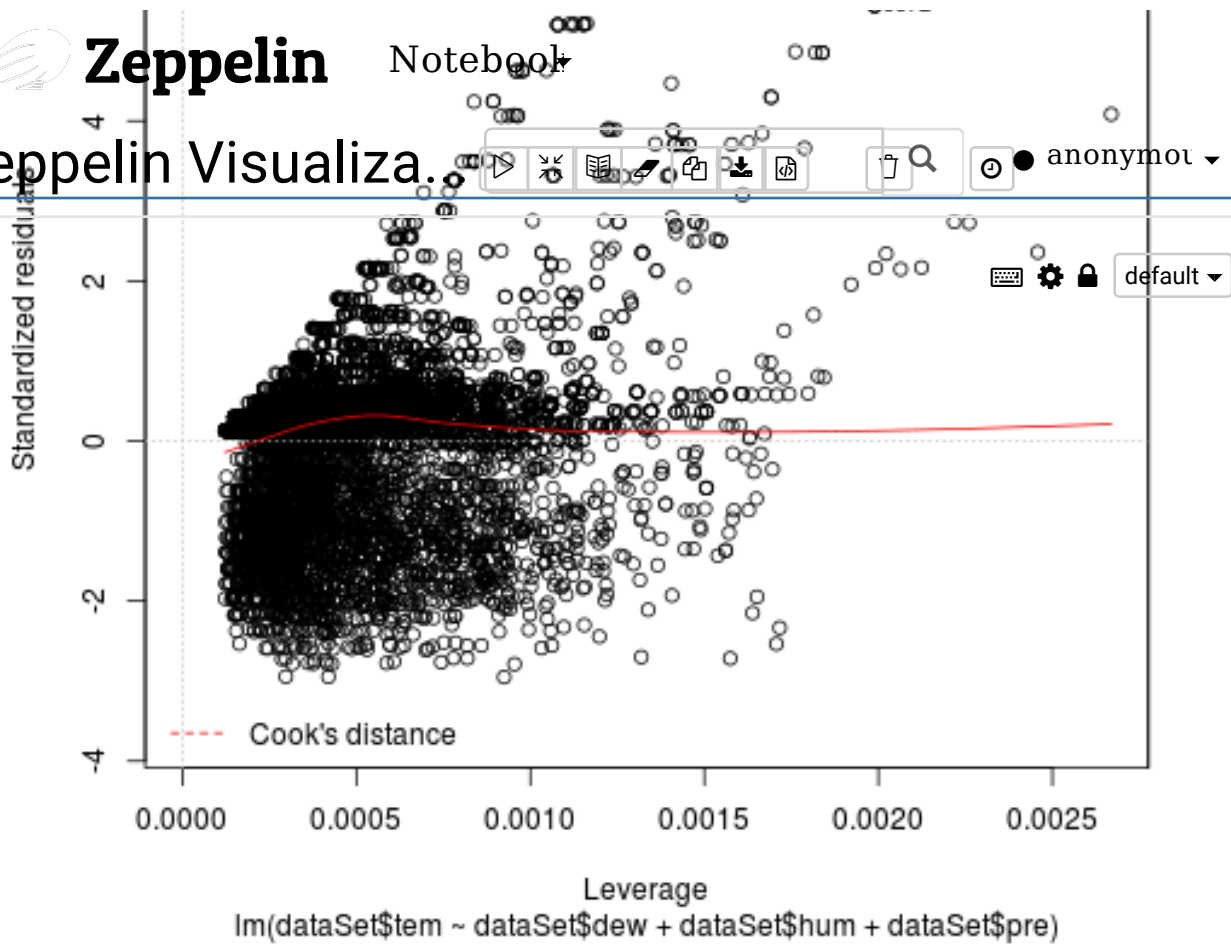




Zeppelin

Notebook

Zeppelin Visualiza.

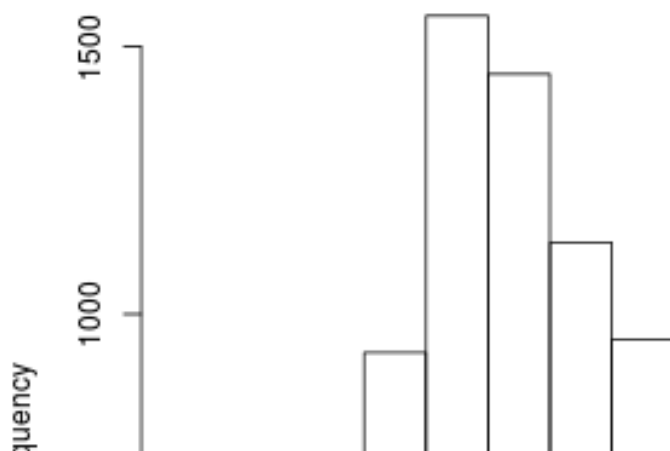


Took 1 sec. Last updated by anonymous at April 18 2017, 7:17:10 PM.

```
%spark.r
hist(dataSet$stem)
```

FINISHED ▶ ✖ 📖 ⚙

Histogram of dataSet\$stem

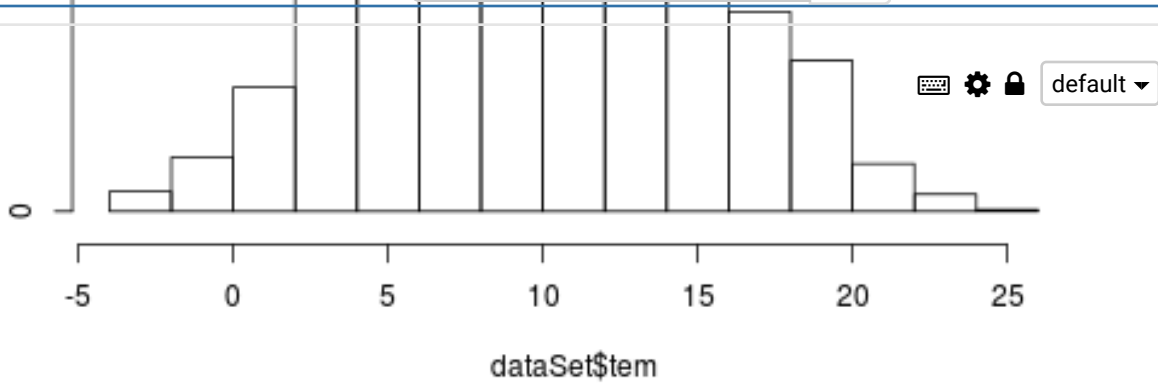




Zeppelin

Notebook

Zeppelin Visualiza...



Took 0 sec. Last updated by anonymous at April 18 2017, 7:28:06 PM.

```
%spark.r
install.packages("corrplot", repos = "http://cran.us.r-project.org")
```

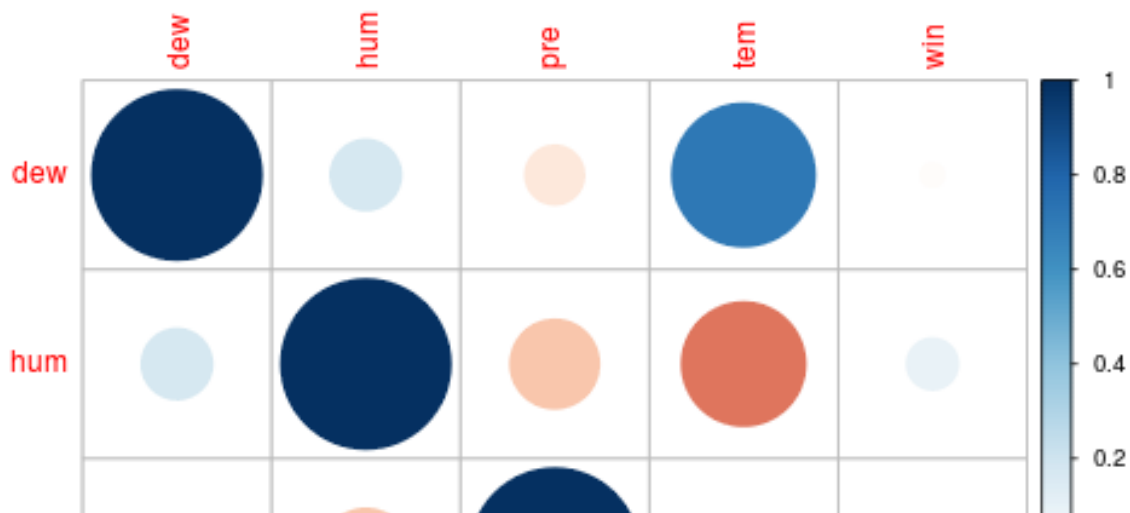
FINISHED ▶ 🔍 📖 ⚙️

The downloaded source packages are in
'/tmp/RtmpWXXF7Y/downloaded_packages'

Took 8 sec. Last updated by anonymous at April 18 2017, 7:44:34 PM.

```
%spark.r
library(corrplot)
M <- cor(dataSet[3:7])
corrplot(M, method="circle")
```

FINISHED ▶ 🔍 📖 ⚙️

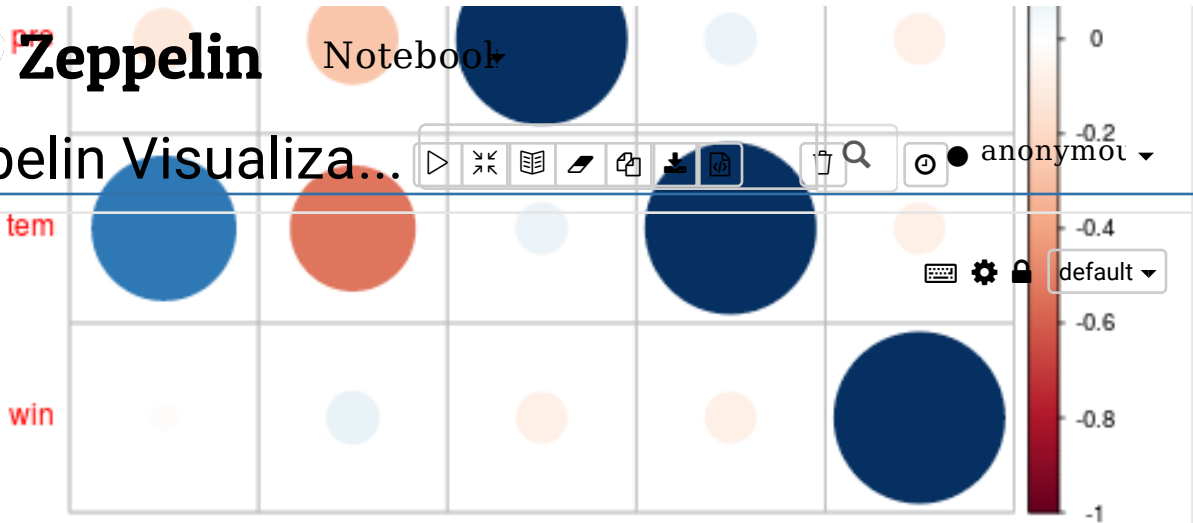




Zeppelin

Notebook

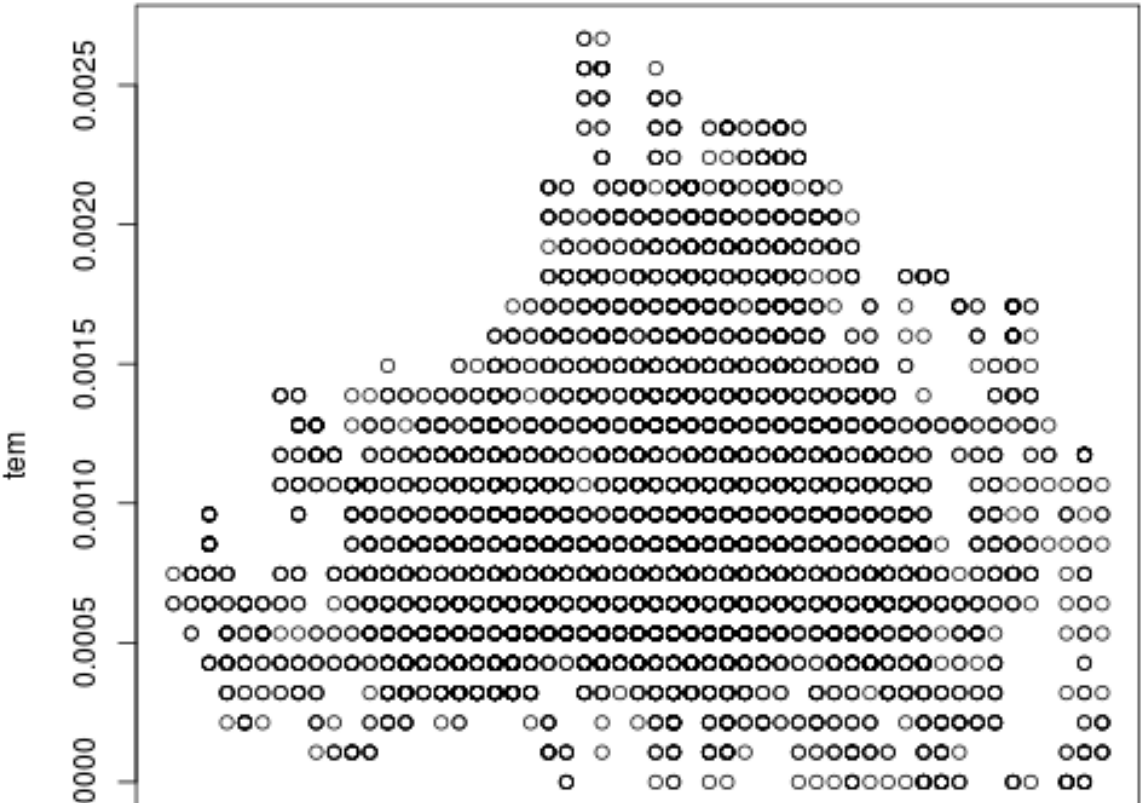
Zeppelin Visualiza...



Took 0 sec. Last updated by anonymous at April 18 2017, 7:45:25 PM.

%spark.r

FINISHED ▶ 🔍 📖 ⚙️



pre

⌨ ⚙ 🔒 default ▾

Took 1 sec. Last updated by anonymous at April 18 2017, 7:49:37 PM. (outdated)

READY    