

# Projeto Data Mining I

*Lucas Parada, Ana Catarina Monteiro, Lucas de Paula*

*11/16/2019*

## Data importation, clean-up and pre-processing

### Importando os Dados

Inicialmente importamos todos as bibliotecas que utilizaremos neste trabalho. Foi importado o dataset “PRSA\_Data\_Aotizhongxin\_20130301-20170228” - Foi utilizado o DataFrame do R para manipular os dados, pois este tipo de estrutura de dados possui um conjunto de funcionalidades e ferramentas que auxiliam neste processo.

### Analizando os dados

Utilizando a função “summary” da linguagem R foi feita uma análise dos dados.

```
##      No          year        month        day
##  Min.   : 1   Min.   :2013   Min.   : 1.000   Min.   : 1.00
##  1st Qu.: 8767  1st Qu.:2014   1st Qu.: 4.000   1st Qu.: 8.00
##  Median :17532  Median :2015   Median : 7.000   Median :16.00
##  Mean   :17532  Mean   :2015   Mean   : 6.523   Mean   :15.73
##  3rd Qu.:26298  3rd Qu.:2016   3rd Qu.:10.000   3rd Qu.:23.00
##  Max.   :35064  Max.   :2017   Max.   :12.000   Max.   :31.00
##
##      hour        PM2.5        PM10        SO2
##  Min.   : 0.00  Min.   : 3.00  Min.   : 2.0  Min.   : 0.2856
##  1st Qu.: 5.75  1st Qu.:22.00  1st Qu.: 38.0  1st Qu.: 3.0000
##  Median :11.50  Median : 58.00  Median : 87.0  Median : 9.0000
##  Mean   :11.50  Mean   : 82.77  Mean   :110.1  Mean   : 17.3759
##  3rd Qu.:17.25  3rd Qu.:114.00 3rd Qu.:155.0  3rd Qu.: 21.0000
##  Max.   :23.00  Max.   :898.00  Max.   :984.0  Max.   :341.0000
##      NA's   :925    NA's   :718    NA's   :935
##      NO2          CO          O3          TEMP
##  Min.   : 2.00  Min.   : 100  Min.   : 0.2142  Min.   :-16.80
##  1st Qu.: 30.00 1st Qu.: 500  1st Qu.: 8.0000  1st Qu.: 3.10
##  Median : 53.00 Median : 900  Median : 42.0000  Median : 14.50
##  Mean   : 59.31 Mean   :1263  Mean   : 56.3534  Mean   : 13.58
##  3rd Qu.: 82.00 3rd Qu.:1500 3rd Qu.: 82.0000  3rd Qu.: 23.30
##  Max.   :290.00 Max.   :10000 Max.   :423.0000  Max.   : 40.50
##  NA's   :1023  NA's   :1776  NA's   :1719  NA's   :20
##      PRES         DEWP        RAIN        wd
##  Min.   : 985.9  Min.   :-35.300  Min.   : 0.00000  NE   : 5140
##  1st Qu.:1003.3  1st Qu.: -8.100  1st Qu.: 0.00000  ENE  : 3950
##  Median :1011.4  Median :  3.800  Median : 0.00000  SW   : 3359
##  Mean   :1011.8  Mean   :  3.123  Mean   : 0.06742  E    : 2608
##  3rd Qu.:1020.1  3rd Qu.: 15.600 3rd Qu.: 0.00000  NNE  : 2445
##  Max.   :1042.0  Max.   : 28.500  Max.   :72.50000  (Other):17481
##  NA's   :20       NA's   :20       NA's   :20       NA's   : 81
##      WSPM        station
##  Min.   : 0.000  Aotizhongxin:35064
##  1st Qu.: 0.900
```

```

## Median : 1.400
## Mean   : 1.708
## 3rd Qu.: 2.200
## Max.   :11.200
## NA's    :14

```

Começamos por verificar se existia algum dia em falta no dataframe e vimos que não. Samendo que o dataset possui os valores referentes a 4 anos completos especificados por hora então devem existir  $(4365+1)24 = 35064$  rows

```

#sabendo que o dataset possui os valores referentes a 4 anos completos especificados por hora então se
dim(df)

```

```

## [1] 35064     18

```

## Outliers

Seguimos com a análise da existência de outliers em variáveis numéricas. Começando por ver fazer a análise por variável como um todo, em seguida fizemos a análise por variável tendo em conta a estação do ano e por último por variável por mês.

### Como um todo

```

col_list = c("CO", "PM2.5", "PM10", "SO2", "NO2", "CO", "O3", "TEMP", "PRES", "DEWP", "RAIN", "WSPM")

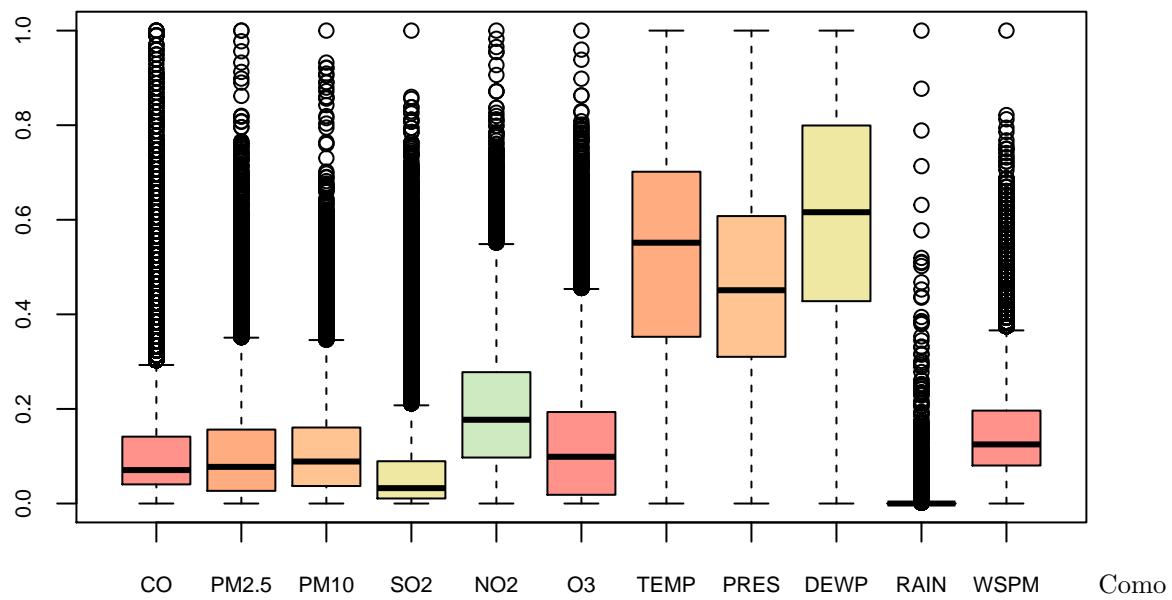
df_ALL_boxplot <- df %>% select(col_list)

# remove rows with missing values
df_ALL_boxplot <- df_ALL_boxplot[ complete.cases(df_ALL_boxplot), ]

preproc<-preProcess(df_ALL_boxplot, method=c("range"))
df_ALL_boxplot <- predict(preproc, df_ALL_boxplot)

boxplot(x = df_ALL_boxplot, col = color_list, cex.axis=0.7)

```



## Por estação do ano

```
season <- function(month, day){  
  print(month)  
  if(month == 12 & day >= 21)  
    return('Winter')  
  if(month == 1 | month == 2)  
    return('Winter')  
  if(month == 3 & day < 20)  
    return('Winter')  
  if(month == 3 & day >= 20)  
    return('Spring')  
  if(month == 4 | month == 5)  
    return('Spring')  
  if(month == 6 & day < 21)  
    return('Spring')  
  if(month == 6 & day >= 21)  
    return('Summer')  
  if(month == 7 | month == 8)  
    return('Summer')  
  if(month == 9 & day < 21)  
    return('Summer')  
  if(month == 9 & day >= 21)  
    return('Autumn')  
  if(month == 10 | month == 11)  
    return('Autumn')  
  if(month == 12 & day < 2)  
    return('Autumn')  
}  
  
col_list <- c("season", "CO", "PM2.5", "PM10", "SO2", "NO2", "CO", "O3", "TEMP", "PRES", "DEWP", "RAIN",  
df_seasons_boxplot <- df %>% mutate(season =  
  ifelse(month == 12 & day >= 21, 'Winter',  
  ifelse(month == 1 | month == 2, 'Winter',  
  ifelse(month == 3 & day < 20, 'Winter',  
  
  ifelse(month == 3 & day >= 20, 'Spring',  
  ifelse(month == 4 | month == 5, 'Spring',  
  ifelse(month == 6 & day < 21, 'Spring',  
  
  ifelse(month == 6 & day >= 21, 'Summer',  
  ifelse(month == 7 | month == 8, 'Summer',  
  ifelse(month == 9 & day < 21, 'Summer',  
  
  ifelse(month == 9 & day >= 21, 'Autumn',  
  ifelse(month == 10 | month == 11, 'Autumn',  
  ifelse(month == 12 & day < 21, 'Autumn',  
  0)))))))))))  
  ) %>%  
  select(col_list)  
  
for(s in df_seasons_boxplot$season %>% unique()){  
  df_aux <- df_seasons_boxplot %>% filter(season == s) %>% select(-season)
```

```

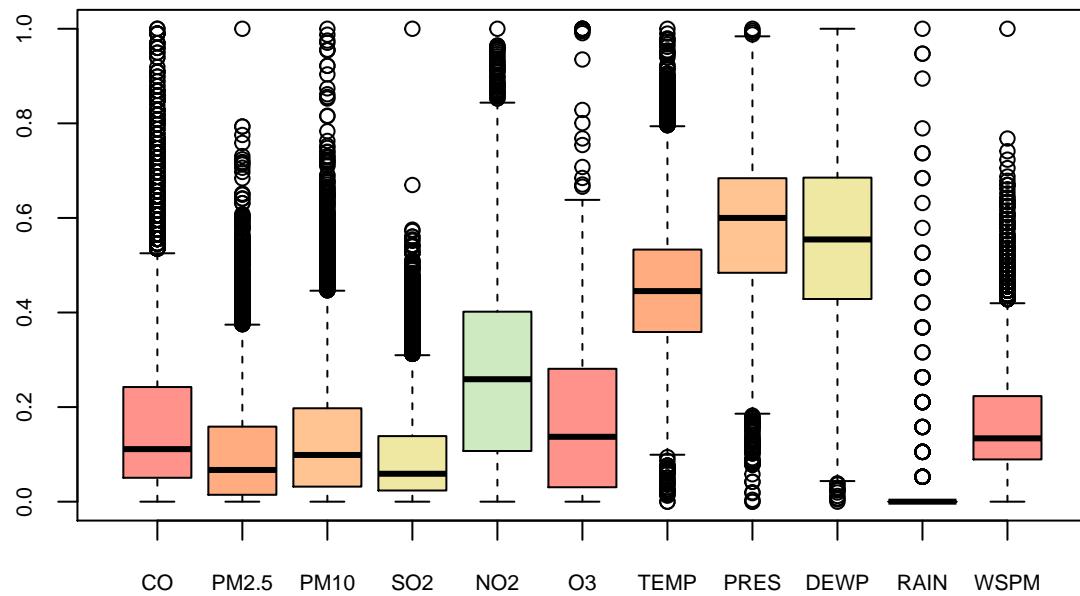
preproc<-preProcess(df_aux, method=c("range"))
norm2 <- predict(preproc, df_aux)

norm2 %>% boxplot(main=s, col = color_list, cex.axis=0.7)
print('')

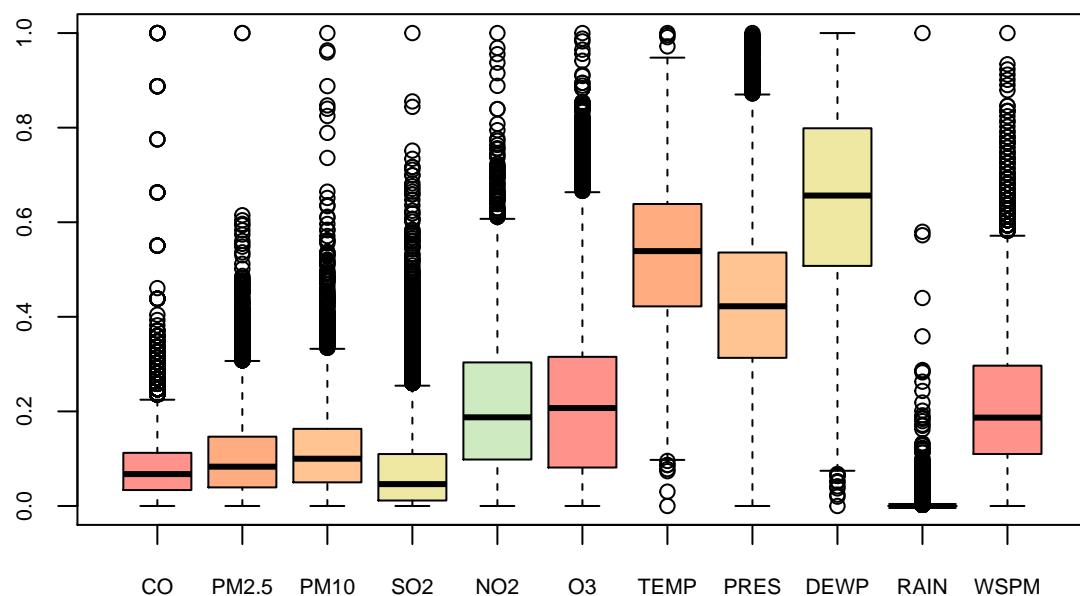
#df_seasons_boxplot %>% filter(season == s) %>% select(-season) %>% scale() %>% boxplot(main=s)
}

```

## Winter

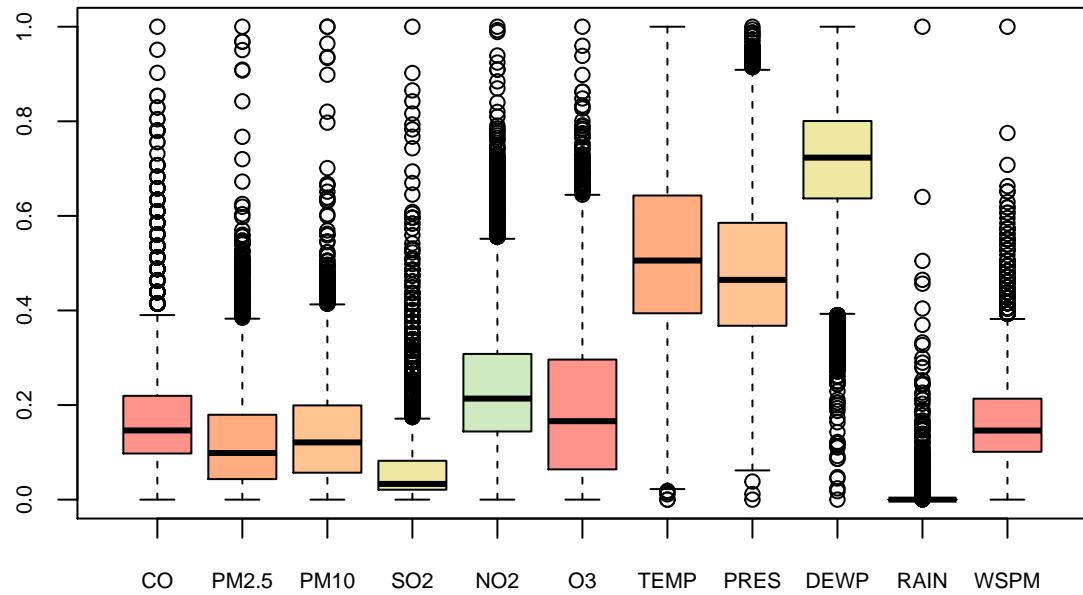


## Spring



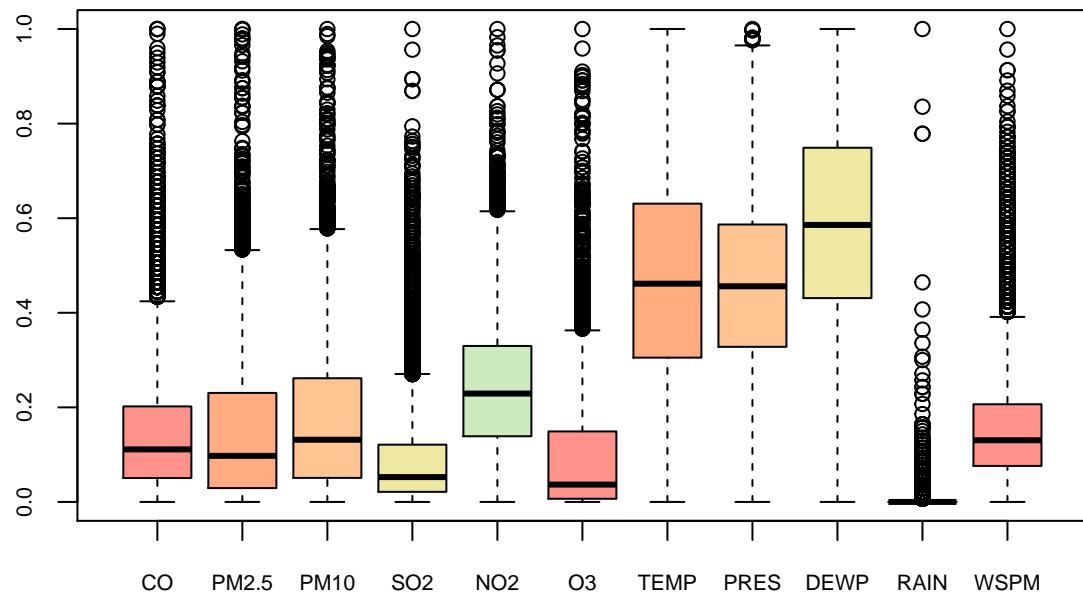
```
## [1] "
```

## Summer



```
## [1] "
```

## Autumn



```
## [1] "
```

Por mês

```
col_list = c("month", "CO", "PM2.5", "PM10", "SO2", "NO2", "CO", "O3", "TEMP", "PRES", "DEWP", "RAIN", "WSPM")  
df_month_boxplot <- df %>% select(col_list)
```

```

month_names <- c("January", "February", "March", "April", "May", "June", "July", "August", "September", "October", "November", "December")

for(m in df_month_boxplot$month %>% unique() %>% sort()){
  df_aux <- df_month_boxplot %>% filter(month == m) %>% select(-month)

  preproc<-preProcess(df_aux, method=c("range"))
  norm2 <- predict(preproc, df_aux)

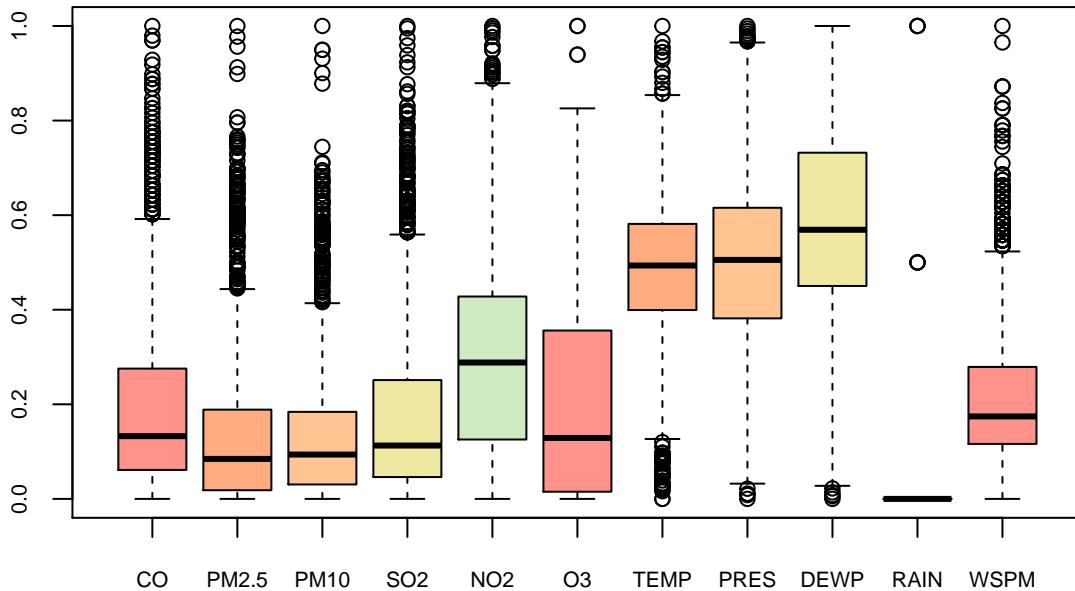
  par(cxy = c(.5,.5))

  norm2 %>% boxplot(main=month_names[m], col = color_list, cex.axis=0.7)
  print('\n')
  #df_month_boxplot %>% filter(month == m) %>% select(-month) %>% scale() %>% boxplot(main=month_names
}

## Warning in par(cxy = c(0.5, 0.5)): graphical parameter "cxy" cannot be set
## [1] "\n"
## Warning in par(cxy = c(0.5, 0.5)): graphical parameter "cxy" cannot be set

```

## January

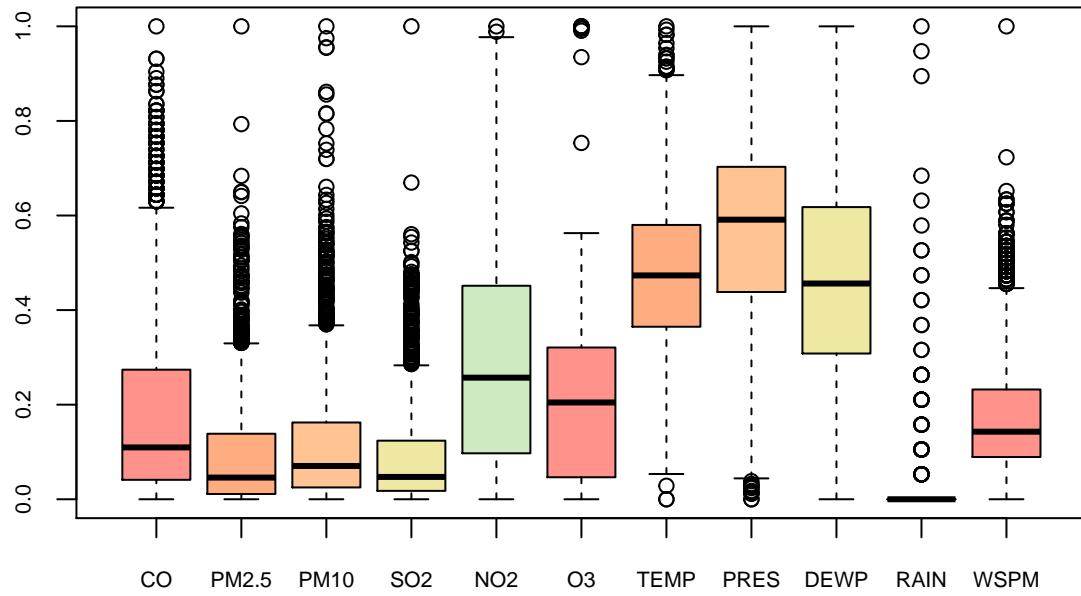


```

## [1] "\n"
## Warning in par(cxy = c(0.5, 0.5)): graphical parameter "cxy" cannot be set

```

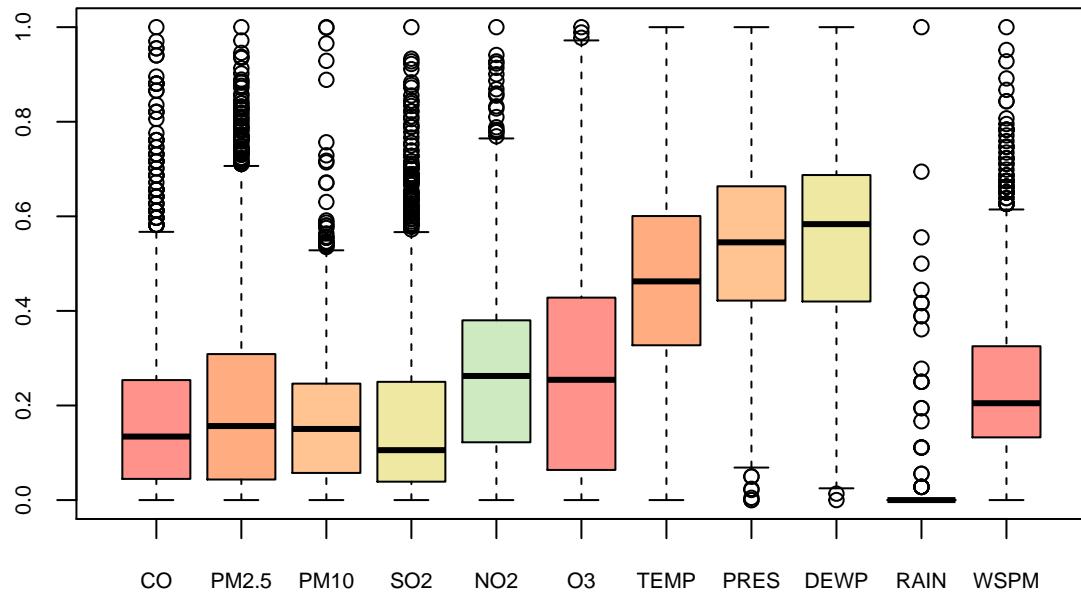
## February



```
## [1] "\n"
```

```
## Warning in par(cxy = c(0.5, 0.5)): graphical parameter "cxy" cannot be set
```

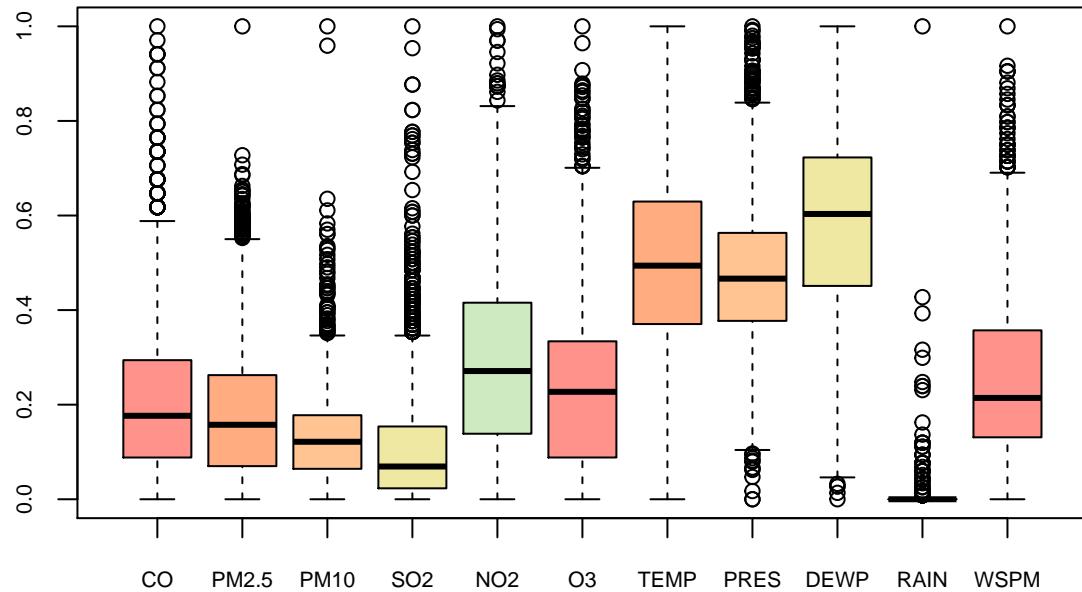
## March



```
## [1] "\n"
```

```
## Warning in par(cxy = c(0.5, 0.5)): graphical parameter "cxy" cannot be set
```

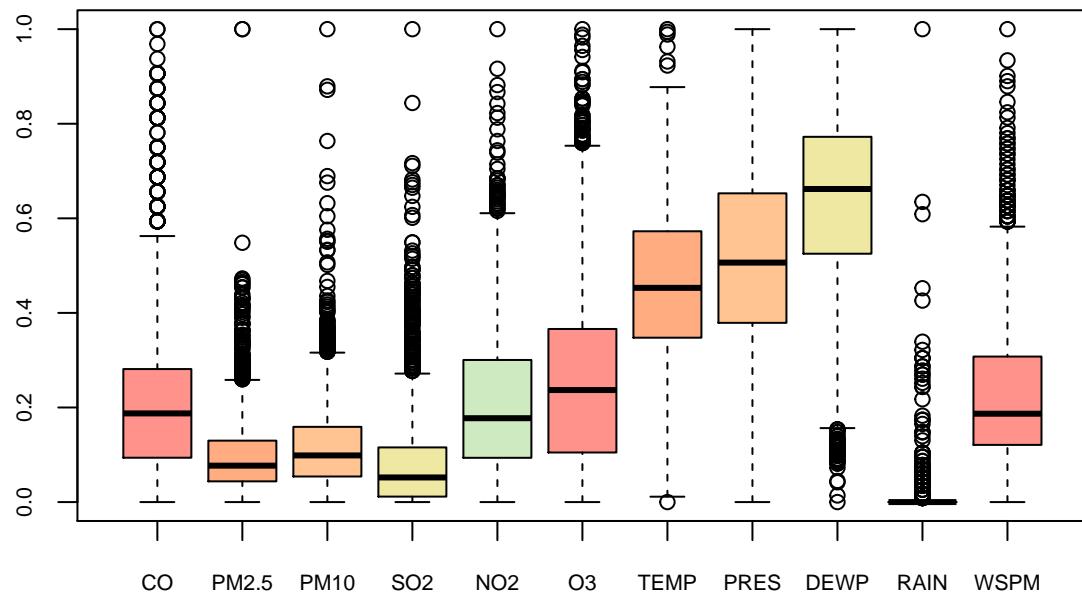
## April



```
## [1] "\n"
```

```
## Warning in par(cxy = c(0.5, 0.5)): graphical parameter "cxy" cannot be set
```

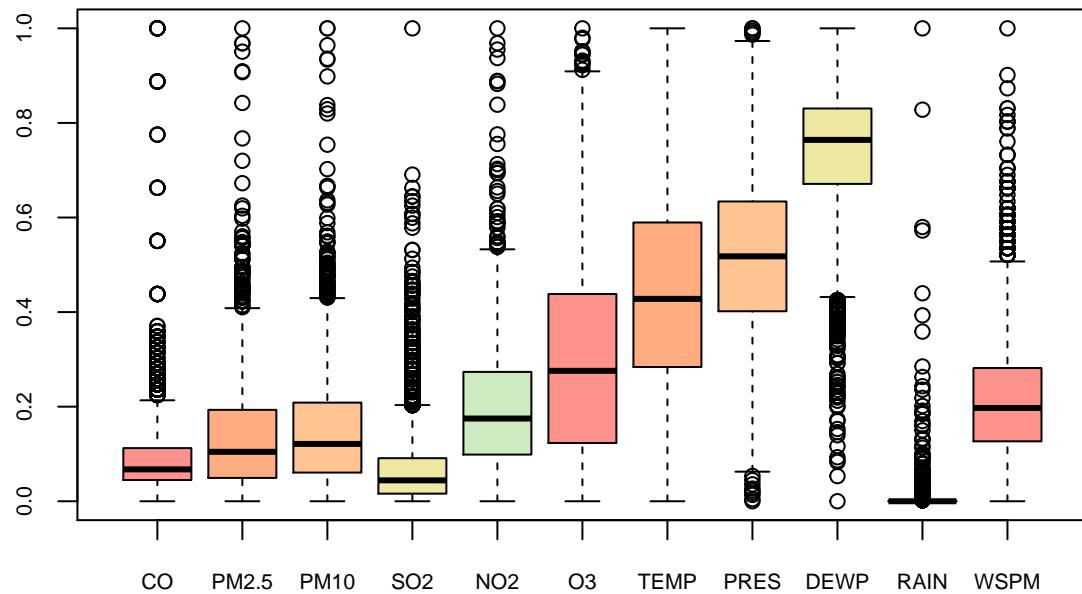
## May



```
## [1] "\n"
```

```
## Warning in par(cxy = c(0.5, 0.5)): graphical parameter "cxy" cannot be set
```

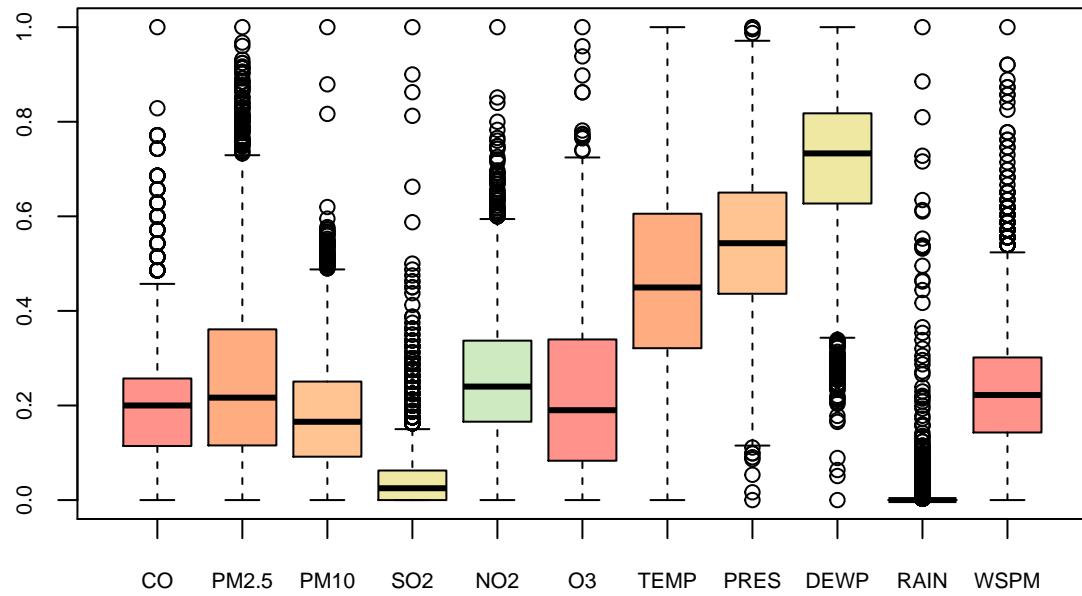
## June



```
## [1] "\n"
```

```
## Warning in par(cxy = c(0.5, 0.5)): graphical parameter "cxy" cannot be set
```

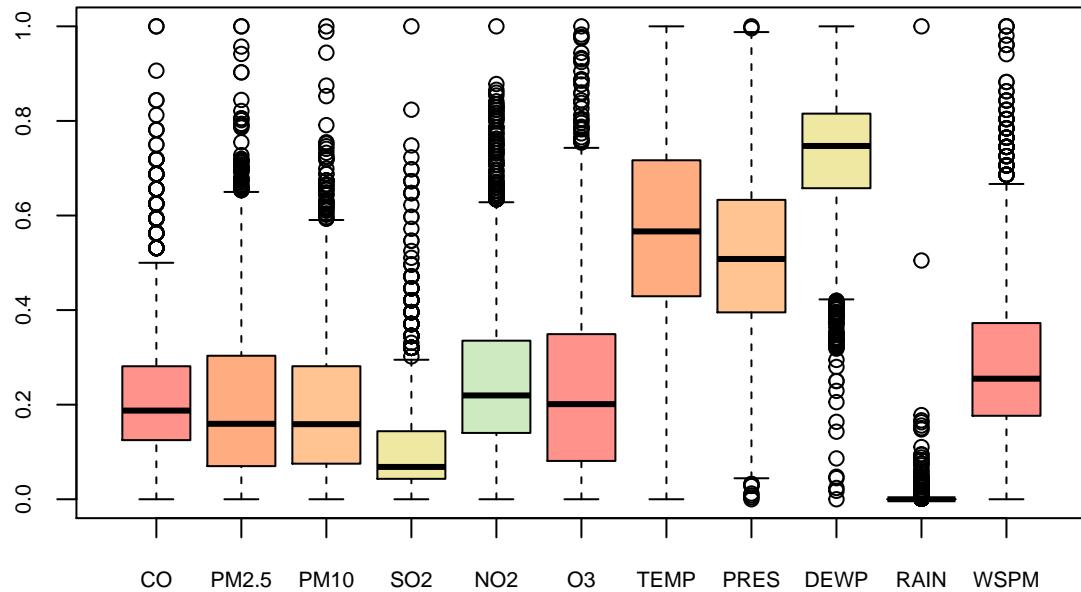
## July



```
## [1] "\n"
```

```
## Warning in par(cxy = c(0.5, 0.5)): graphical parameter "cxy" cannot be set
```

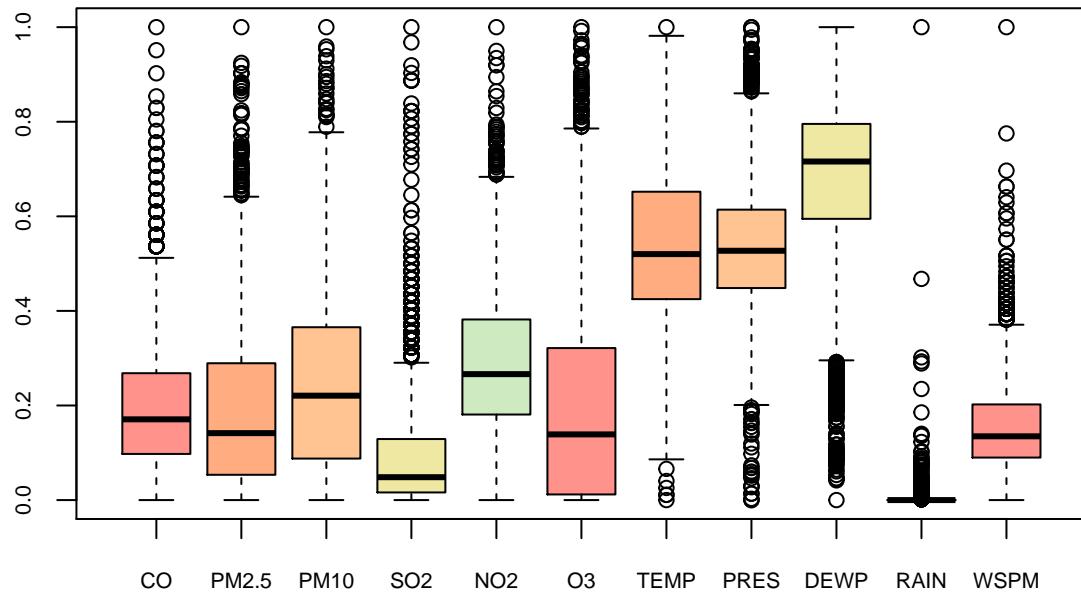
## August



```
## [1] "\n"
```

```
## Warning in par(cxy = c(0.5, 0.5)): graphical parameter "cxy" cannot be set
```

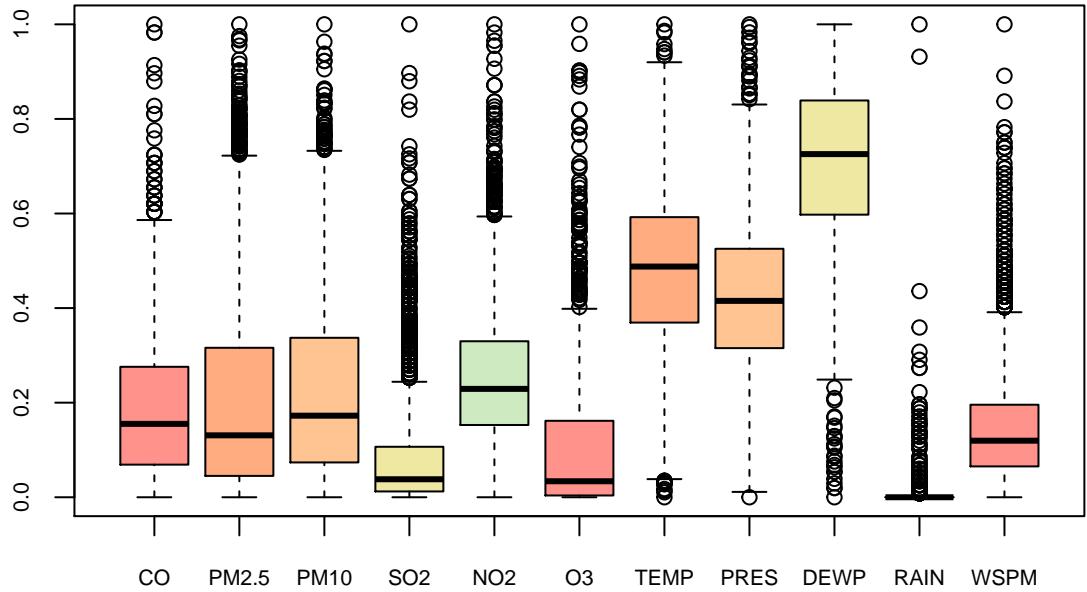
## September



```
## [1] "\n"
```

```
## Warning in par(cxy = c(0.5, 0.5)): graphical parameter "cxy" cannot be set
```

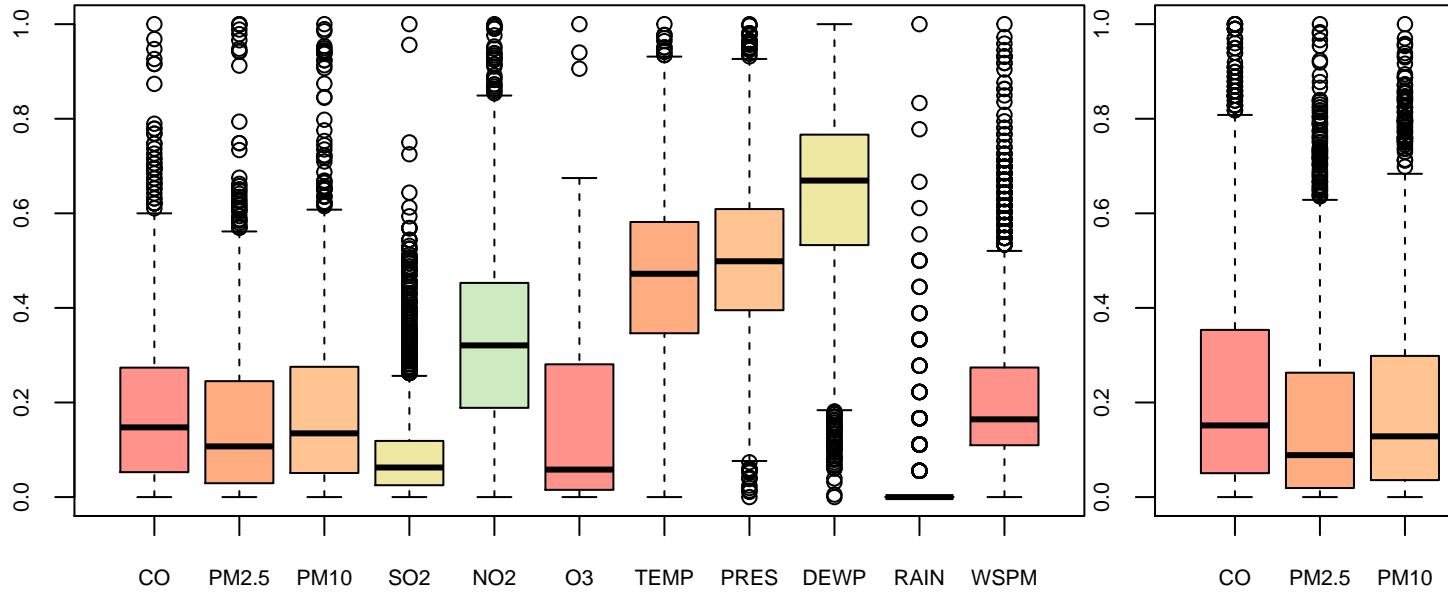
## October



```
## [1] "\n"
```

```
## Warning in par(cxy = c(0.5, 0.5)): graphical parameter "cxy" cannot be set
```

## November



```
## [1] "\n"
```

```
#df_month_boxplot %>% group_by(month) %>% select(-month) %>% scale() %>% boxplot()
```

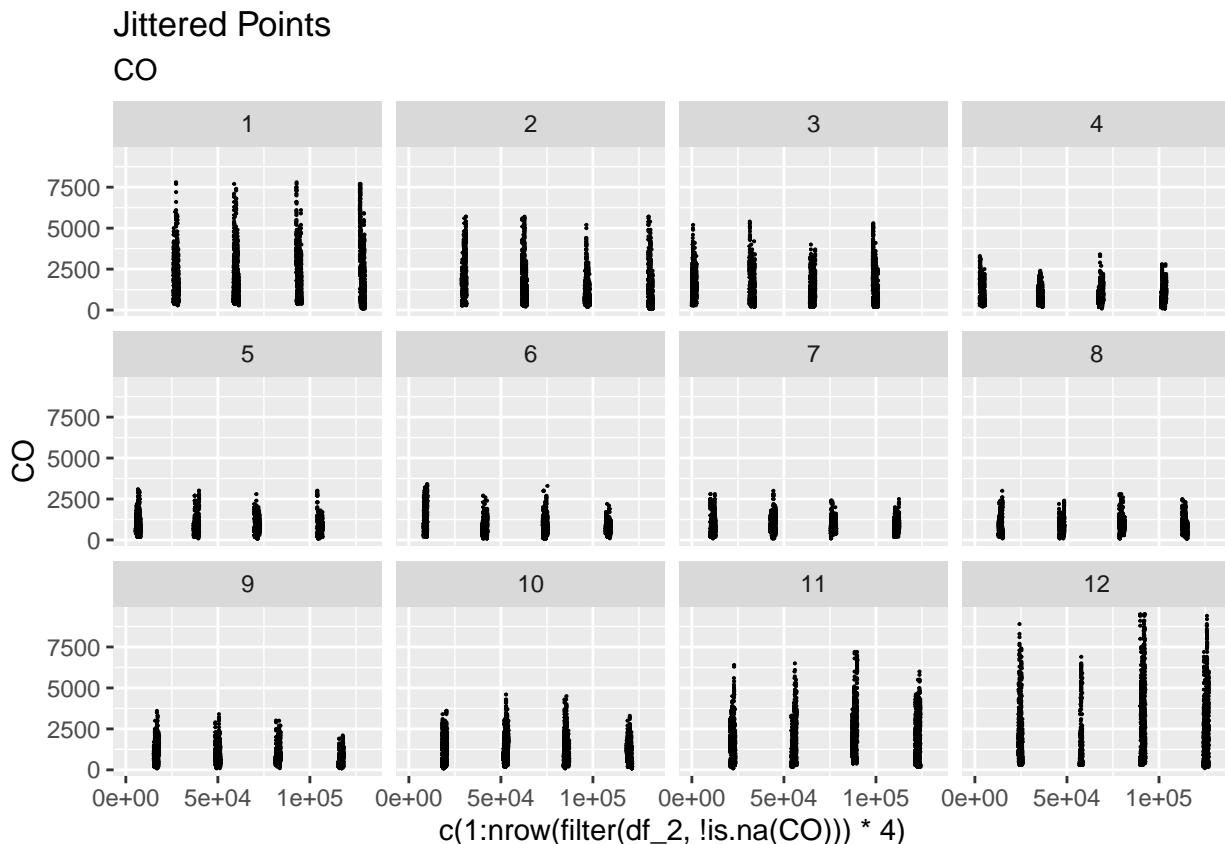
```
#df %>% group_by(month) %>% ggplot(aes(x = month, y = TEMP)) + geom_boxplot() + facet_wrap(~year)
```

## Tratamento de outliers

```
#identificar se é melhor retirar os outliers como um todo, por estação do ano ou por mês
col_list = c("CO", "PM2.5", "PM10", "SO2", "NO2", "O3", "TEMP", "PRES", "DEWP", "RAIN", "WSPM")
df_2 <- df

df_2 <- df_2 %>% group_by(month) %>%
  mutate(CO = ifelse(CO > 4*mean(CO, na.rm = T), NA, CO)) %>%
  mutate(PM2.5 = ifelse(PM2.5 > 4*mean(PM2.5, na.rm = T), NA, PM2.5)) %>%
  mutate(PM10 = ifelse(PM10 > 4*mean(PM10, na.rm = T), NA, PM10)) %>%
  mutate(SO2 = ifelse(SO2 > 4*mean(SO2, na.rm = T), NA, SO2)) %>%
  mutate(NO2 = ifelse(NO2 > 4*mean(NO2, na.rm = T), NA, NO2)) %>%
  mutate(O3 = ifelse(O3 > 4*mean(O3, na.rm = T), NA, O3)) %>%
  mutate(TEMP = ifelse(TEMP > 4*mean(TEMP, na.rm = T), NA, TEMP)) %>%
  mutate(PRES = ifelse(PRES > 4*mean(PRES, na.rm = T), NA, PRES)) %>%
  mutate(DEWP = ifelse(DEWP > 4*mean(DEWP, na.rm = T), NA, DEWP)) %>%
#  mutate(RAIN = ifelse(RAIN > 4*mean(RAIN, na.rm = T), NA, RAIN)) %>%
#  mutate(WSPM = ifelse(WSPM > 4*mean(WSPM, na.rm = T), NA, WSPM)) %>%
ungroup()

df_2 %>% filter(!is.na(CO)) %>% ggplot(aes(c(1:nrow(filter(df_2, !is.na(CO)))*4), CO)) + geom_point(siz
```



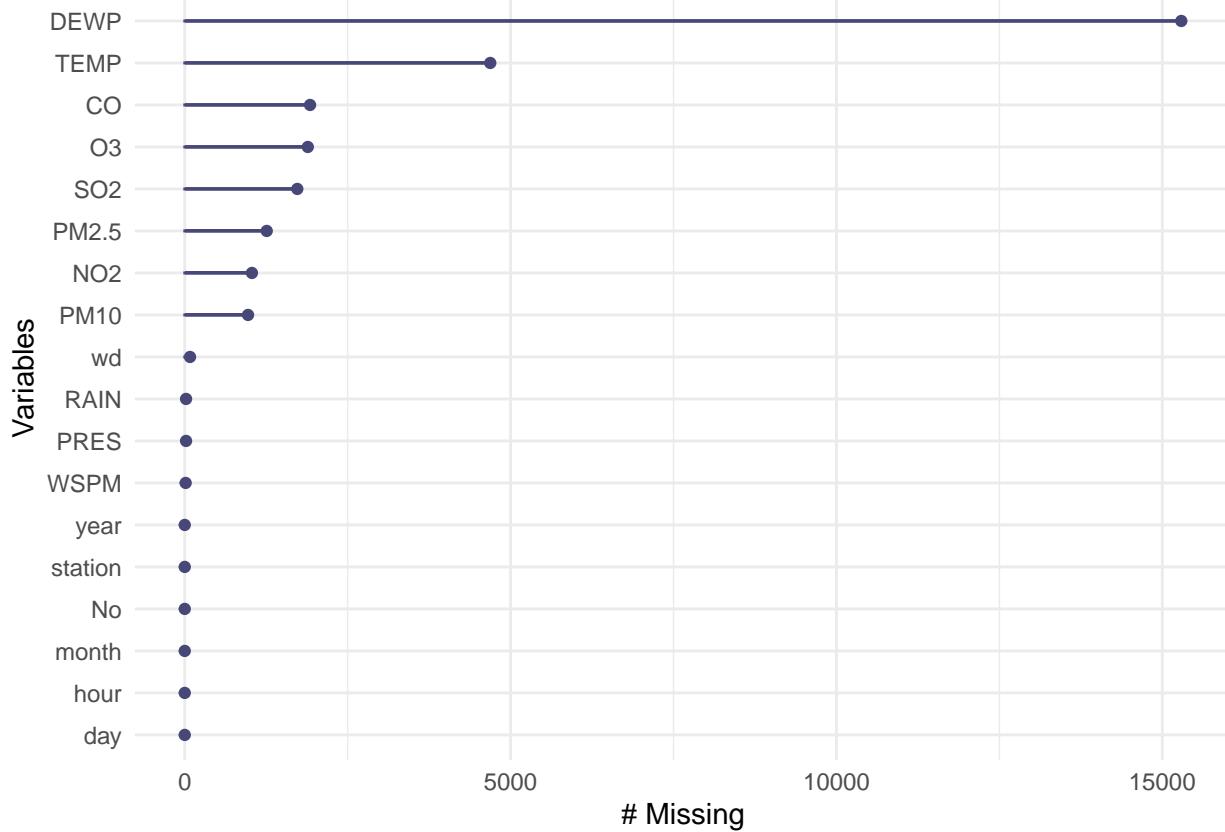
## One

Apartir da análise mencionada anteriormente foi identificada que as varáveis “PM2.5”, “PM10”, “SO2”, “NO2”, “CO”, “O3”, “TEMP”, “PRES”, “DEWP”, “RAIN”, “wd”, “WSPM” possuem missing values.

Em seguida foi uma análise e foi verificado que as variáveis “TEMP”, “PRES”, “DEWP”, “RAIN” esta a faltar no mesmo dia através do gráfico:

Como a variável DEWP possui muitos missings values fizemos um estudo relacionando a varável DEWP e as outras variáveis climáticas para assim tentarmos identificar alguma relação entre elas. Concluimos que existe uma relação direta com a TEMP e uma relação inversa com A variável PRES.

```
#check missing values
#df %>% select(RAIN, TEMP, PRES, wd, WSPM) %>% filter_any_na()
gg_miss_var(df_2)
```



```
atm = c("TEMP", "PRES", "DEWP", "RAIN")
# TEMP
theme_set(theme_bw()) # pre-set the bw theme.
ggplot(df_2, aes(TEMP, DEWP)) + geom_jitter(width = .5, size=.1) +
  labs(
    title="TEMP",
    subtitle="Relação entre TEMP e DEWP",
    x="DEWP",
    y="TEMP") +
  geom_count() +
  geom_smooth(method="lm", se=F)
```

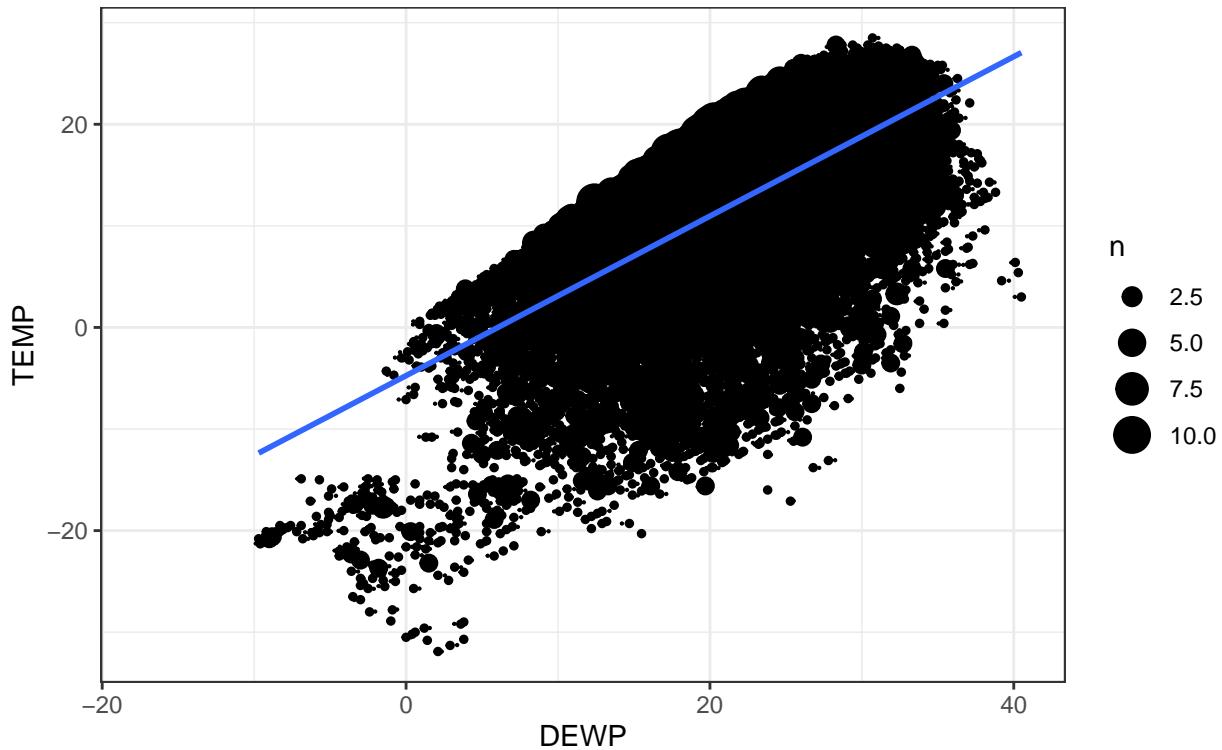
```

## Warning: Removed 15291 rows containing non-finite values (stat_sum).
## Warning: Removed 15291 rows containing non-finite values (stat_smooth).
## Warning: Removed 15291 rows containing missing values (geom_point).

```

## TEMP

Relação entre TEMP e DEWP



```

# PRES
ggplot(df_2, aes(PRES, DEWP)) + geom_jitter(width = .5, size=.1) +
  labs(
    title="PRES",
    subtitle="Relação entre PRES e DEWP",
    x="PRES",
    y="DEWP")+
  geom_count() +
  geom_smooth(method="lm", se=F)

```

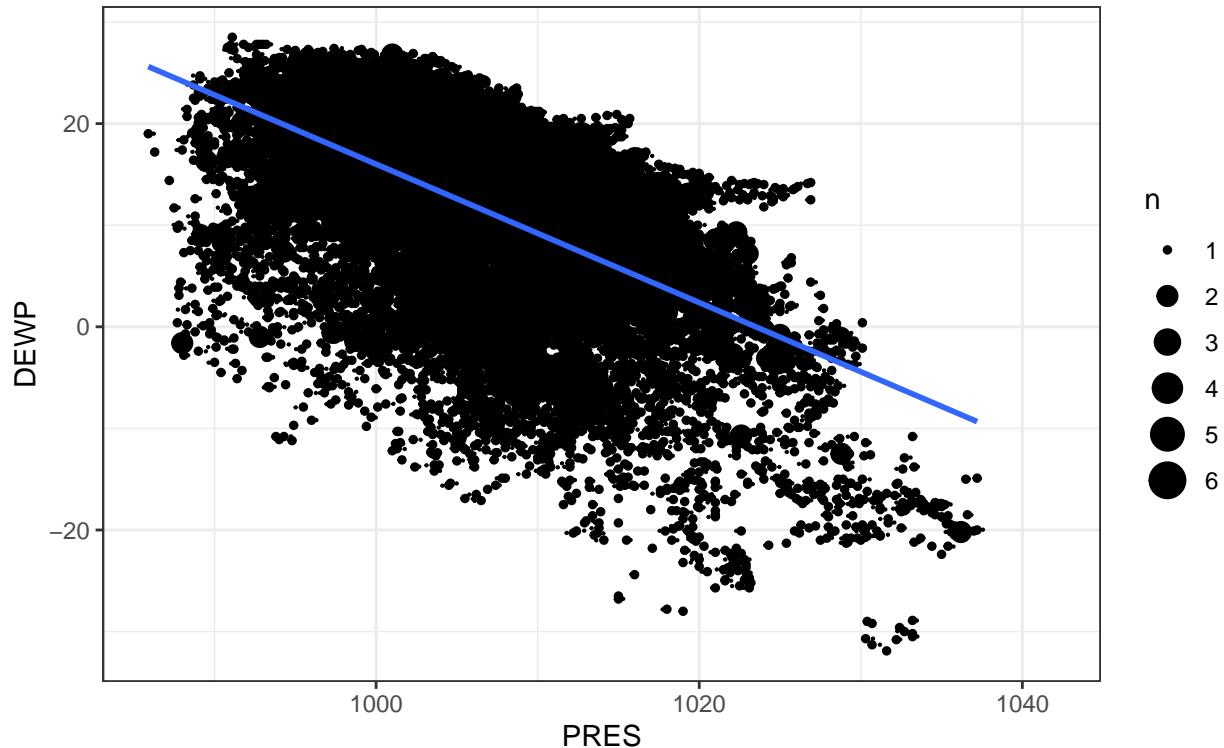
```

## Warning: Removed 15291 rows containing non-finite values (stat_sum).
## Warning: Removed 15291 rows containing non-finite values (stat_smooth).
## Warning: Removed 15291 rows containing missing values (geom_point).

```

## PRES

Relação entre PRES e DEWP

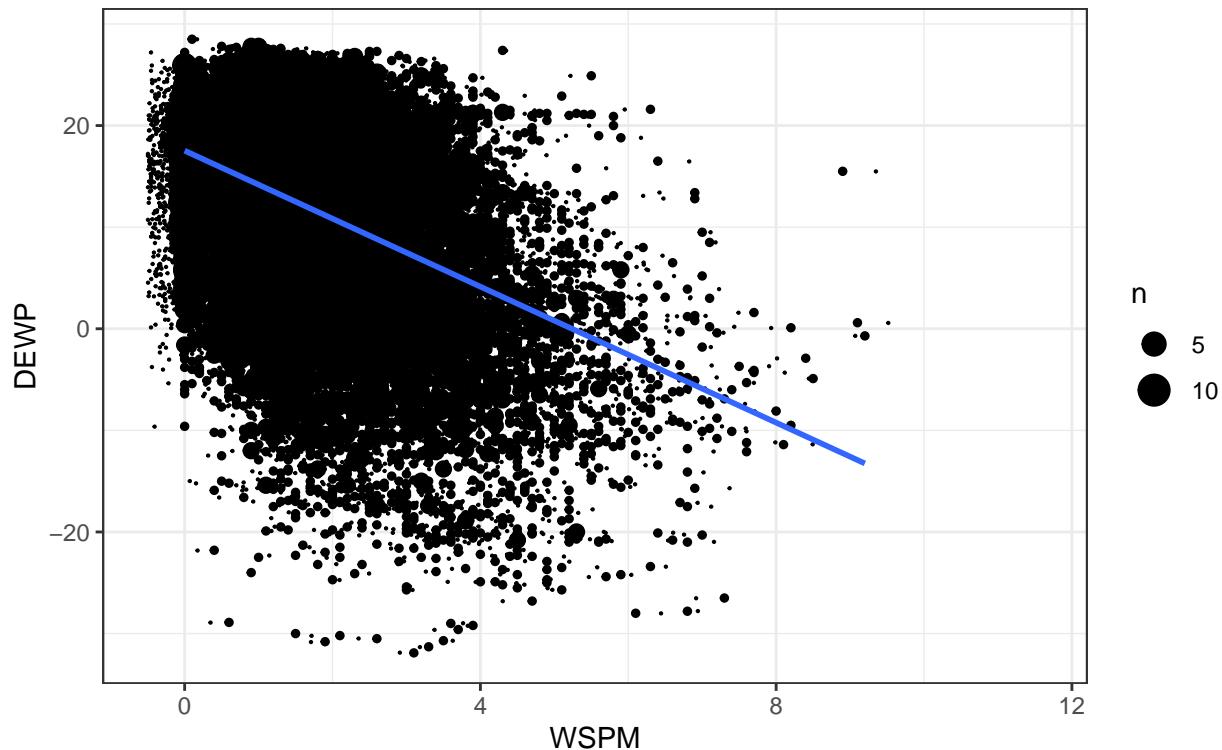


```
# DEWP
ggplot(df_2, aes(WSPM, DEWP)) + geom_jitter(width = .5, size=.1) +
  labs(
    title="Jittered Points",
    subtitle="Relação entre WSPM e DEWP",
    x="WSPM",
    y="DEWP") +
  geom_count() +
  geom_smooth(method="lm", se=F)

## Warning: Removed 15291 rows containing non-finite values (stat_sum).
## Warning: Removed 15291 rows containing non-finite values (stat_smooth).
## Warning: Removed 15291 rows containing missing values (geom_point).
```

## Jittered Points

### Relação entre WSPM e DEWP

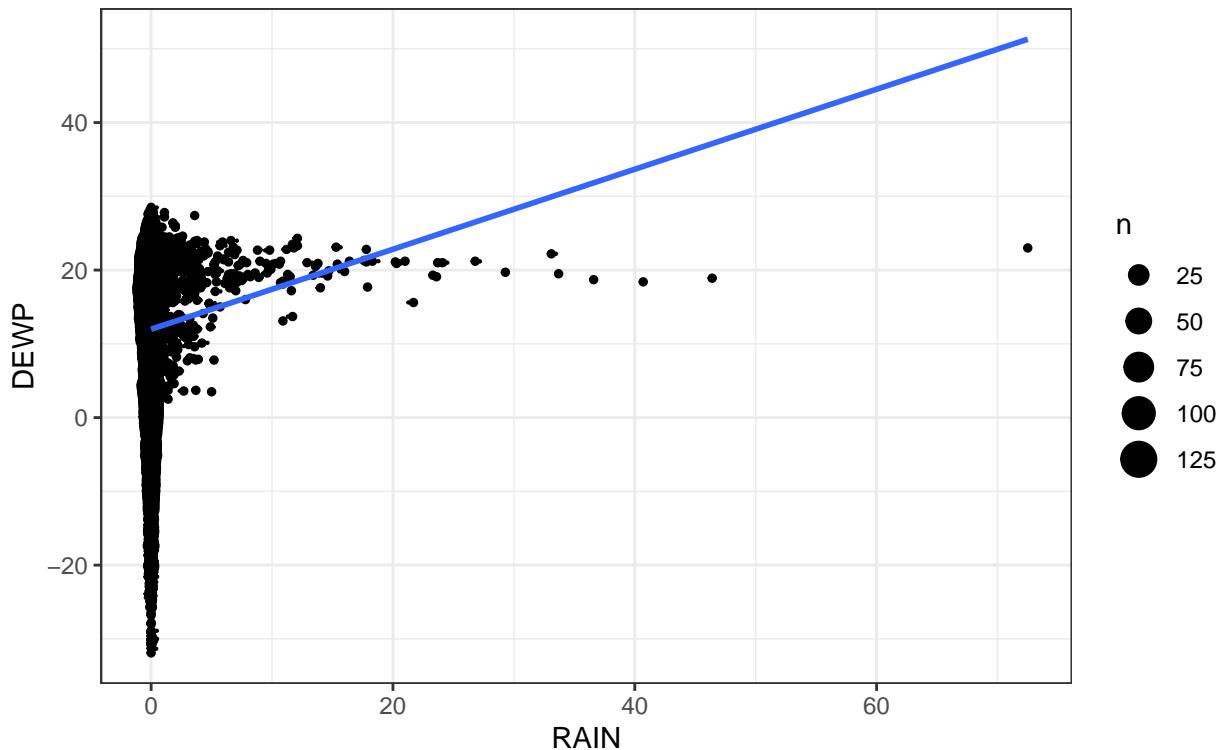


```
# RAIN
ggplot(df_2, aes(RAIN, DEWP)) + geom_jitter(width = .5, size=.1) +
  labs(
    title="RAIN",
    subtitle="Relação entre RAIN e DEWP",
    x="RAIN",
    y="DEWP") +
  geom_count() +
  geom_smooth(method="lm", se=F)

## Warning: Removed 15291 rows containing non-finite values (stat_sum).
## Warning: Removed 15291 rows containing non-finite values (stat_smooth).
## Warning: Removed 15291 rows containing missing values (geom_point).
```

## RAIN

### Relação entre RAIN e DEWP



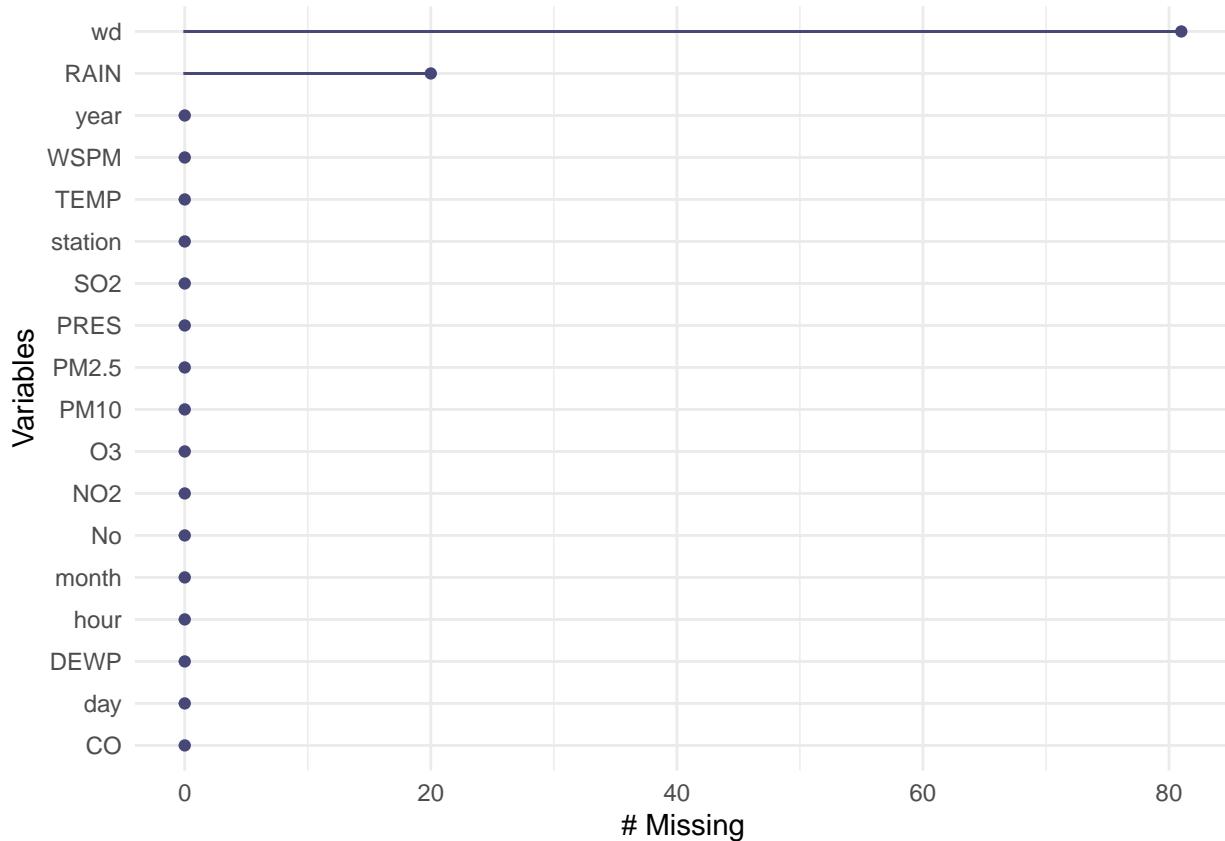
```
#for(col_name in c('TEMP', 'PRES', 'DEWP', 'RAIN')){  
#  (df %>% mutate(test_col = ifelse(is_na(TEMP), 1, 0)))$test_col %>% plot(pch=". ", cex=2, main=col_name  
#  abline(h = 4*mean(c(df[col_name]), na.rm=T), col="red") # add cutoff line  
#}
```

## Data Analisys

```
#head(df, 20)  
#summary(df)  
#dim(df)
```

## Missing Values

```
interpolation_df <- function(d, col_names ){  
  for(col in col_names){  
    d[[col]] <- na.approx(d[[col]], rule=2)  
  }  
  return(d)  
}  
  
col_names = c("PM2.5", "PM10", "SO2", "NO2", "CO", "O3", "TEMP", "PRES", "DEWP", "WSPM")  
df_2 <- interpolation_df(df_2, col_names)  
  
gg_miss_var(df_2)
```



```

#RAIN
df_2$RAIN[which(is.na(df_2$RAIN))] <- mode(df_2$RAIN)
#WD
df_2 <- transform(df_2, wd = na.locf(wd))

#WSPM
#Precisamo verificar como tratar desta variavel(relação entre temperatura e )

df_2 %>% any_na()

## [1] FALSE

```

### Calcular as unidades de medida apropriadas

```

#massa atomica de atomos utilizados
o <- 15.999
n <- 14.007
c <- 12.011
s <- 32.06

#calculando o volume a partir da pressao e da temperatura
volum_func <- function(t, p){
  return(22.41 * ((t+273.15)/273.15) * 1013/p)
}

# conversão de ug/m^3 para ppb
df_2 <- df_2 %>% mutate( CO = CO * volum_func(TEMP, PRES) / (c+o) )

```

```

df_2 <- df_2 %>% mutate( S02 = S02 * volum_func(TEMP, PRES) / (s+o*2) )
df_2 <- df_2 %>% mutate( N02 = N02 * volum_func(TEMP, PRES) / (n+o*2) )
df_2 <- df_2 %>% mutate( O3 = O3 * volum_func(TEMP, PRES) / (o*3) )

# conversão de ppb para ppm
df_2 <- df_2 %>% mutate( CO = CO / 1000 )

```

## calcular a estação do ano

```

# Selecionamos uma sequencia logica para a estação do ano, criando uma escala de de 0 a 3, sendo 0 a ma
df_2 <- df_2 %>% mutate(season =
  ifelse(month == 12 & day >= 21, 1, #outono
  ifelse(month == 1 | month == 2, 1, #outono
  ifelse(month == 3 & day < 20, 1, #outono

  ifelse(month == 3 & day >= 20, 0.33, #primavera
  ifelse(month == 4 | month == 5, 0.33, #primavera
  ifelse(month == 6 & day < 21, 0.33, #primavera

  ifelse(month == 6 & day >= 21, 0, #verao
  ifelse(month == 7 | month == 8, 0, #verao
  ifelse(month == 9 & day < 21, 0, #verao

  ifelse(month == 9 & day >= 21, 0.66, #inverno
  ifelse(month == 10 | month == 11, 0.66, #inverno
  ifelse(month == 12 & day < 21, 0.66, #inverno
  0)))))))))))

```

## One Hot Encode da variavel wd

```

dmy <- dummyVars(~wd, data = df_2)
trsfc <- data.frame(predict(dmy, newdata = df_2))

df_2 <- cbind(df_2, trsf)

```

## Calcular AQI

```

aqi_table <- read.csv("aqi_table.csv", header = T)
calc_aqi <- function(c, c_h, c_l){
  ii = 1
  i_h <- aqi_table$AQI_NUM_high
  i_l <- aqi_table$AQI_NUM_low

  for(i in c(1:7)){
    if(c_h[i]>=c)
      ii = i
      break;
  }
  return( ((i_h[ii] - i_l[ii])/(c_h[ii] - c_l[ii])) * (c - c_l[ii]) + i_l[ii])
}

```

```

df_2 <- df_2 %>% mutate( CO_AQI      = calc_aqi(CO,aqi_table$CO_high, aqi_table$CO_low) )

## Warning in if (c_h[i] >= c) ii = i: the condition has length > 1 and only
## the first element will be used

df_2 <- df_2 %>% mutate( PM2.5_AQI = calc_aqi(PM2.5, aqi_table$PM2.5_high, aqi_table$PM2.5_low) )

## Warning in if (c_h[i] >= c) ii = i: the condition has length > 1 and only
## the first element will be used

df_2 <- df_2 %>% mutate( PM10_AQI   = calc_aqi(PM10, aqi_table$PM10_high, aqi_table$PM10_low) )

## Warning in if (c_h[i] >= c) ii = i: the condition has length > 1 and only
## the first element will be used

df_2 <- df_2 %>% mutate( SO2_AQI    = calc_aqi(SO2,aqi_table$SO2_high, aqi_table$SO2_low) )

## Warning in if (c_h[i] >= c) ii = i: the condition has length > 1 and only
## the first element will be used

df_2 <- df_2 %>% mutate( NO2_AQI    = calc_aqi(NO2,aqi_table$NO2_high, aqi_table$NO2_low) )

## Warning in if (c_h[i] >= c) ii = i: the condition has length > 1 and only
## the first element will be used

df_2 <- df_2 %>% mutate( O3_AQI     = calc_aqi(O3,aqi_table$O3_high, aqi_table$O3_low) )

## Warning in if (c_h[i] >= c) ii = i: the condition has length > 1 and only
## the first element will be used

```

## Calcular a media diaria das variaveis SO2, NO2, PM10, PM2.5

```

df_2 <- df_2 %>% group_by(year, month, day) %>%
  mutate(SO2_AQI_mean = mean(SO2_AQI, na.rm = T)) %>%
  mutate(NO2_AQI_mean = mean(NO2_AQI, na.rm = T)) %>%
  mutate(PM10_AQI_mean = mean(PM10_AQI, na.rm = T)) %>%
  mutate(PM2_AQI_mean = mean(PM2.5_AQI, na.rm = T)) %>%
  ungroup()
df_2

## # A tibble: 35,064 x 45
##       No year month   day hour PM2.5[, "PM2.5"] PM10[, "PM10"] SO2[, "SO2"]
##   <int> <int> <int> <int> <int>      <dbl>        <dbl>      <dbl>
## 1     1   2013     3     1     0          4          4      1.38
## 2     2   2013     3     1     1          8          8      1.38
## 3     3   2013     3     1     2          7          7      1.72
## 4     4   2013     3     1     3          6          6      3.79
## 5     5   2013     3     1     4          3          3      4.12
## 6     6   2013     3     1     5          5          5      6.17
## 7     7   2013     3     1     6          3          3      6.16
## 8     8   2013     3     1     7          3          6      6.52
## 9     9   2013     3     1     8          3          6      5.52
## 10   10   2013     3     1     9          3          8      4.15
## # ... with 35,054 more rows, and 37 more variables: NO2[, "NO2"] <dbl>,
## #   CO[, "CO"] <dbl>, O3[, "O3"] <dbl>, TEMP[, "TEMP"] <dbl>,
## #   PRES[, "PRES"] <dbl>, DEWP[, "DEWP"] <dbl>, RAIN <chr>, wd <fct>,
## #   WSPM[, "WSPM"] <dbl>, station <fct>, season <dbl>, wd.E <dbl>,

```

```

## # wd.ENE <dbl>, wd.ESE <dbl>, wd.N <dbl>, wd.NE <dbl>, wd.NNE <dbl>,
## # wd.NNW <dbl>, wd.NW <dbl>, wd.S <dbl>, wd.SE <dbl>, wd.SSE <dbl>,
## # wd.SSW <dbl>, wd.SW <dbl>, wd.W <dbl>, wd.WNW <dbl>, wd.WSW <dbl>,
## # CO_AQI[,"CO"] <dbl>, PM2.5_AQI[,"PM2.5"] <dbl>,
## # PM10_AQI[,"PM10"] <dbl>, SO2_AQI[,"SO2"] <dbl>, NO2_AQI[,"NO2"] <dbl>,
## # O3_AQI[,"O3"] <dbl>, SO2_AQI_mean <dbl>, NO2_AQI_mean <dbl>,
## # PM10_AQI_mean <dbl>, PM2_AQI.5_mean <dbl>

```

Calc a class variable

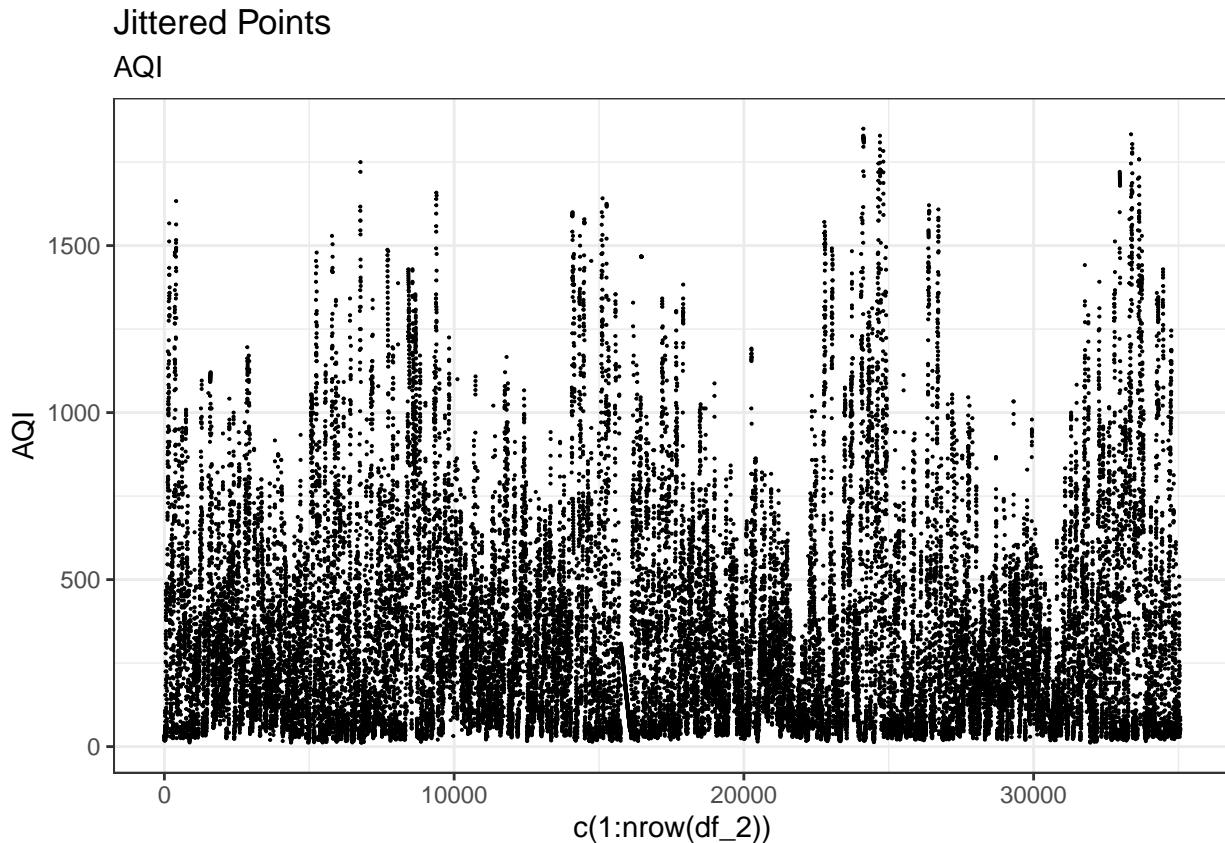
```

df_2[, "AQI"] <- apply( df_2[c("PM2.5_AQI", "PM10_AQI", "SO2_AQI", "NO2_AQI", "CO_AQI", "O3_AQI")], 1, mean)

df_2 <- df_2 %>% mutate(AQI_cat = ifelse(AQI <= 50, 0,
                                             ifelse(AQI <= 100, 1,
                                                   ifelse(AQI <= 150, 2,
                                                       ifelse(AQI <= 200, 3,
                                                         ifelse(AQI <= 300, 4, 5 )))))

df_2 %>% ggplot(aes(c(1:nrow( df_2 )), AQI)) +
  geom_point(size=.05) +
  labs(subtitle="AQI", title="Jittered Points")

```



## Data exploratory analysis

### Predictive modelling

#### installs and imports

```
#install.packages("tensorflow")
#install.packages("keras")

suppressMessages(library(tensorflow))
suppressMessages(library(keras))

tf$constant("Hellow Tensorflow")

## tf.Tensor(Hellow Tensorflow, shape=(), dtype=string)
```