# Principal components of distributional SD

A. Irpino, R. Verde
ESTP Cologne 14-16 May 2024

# Principal Component Analysis of distributional data

We use the `HistDAWass` R package for showing the main procedures.

```
1  # install.packages("HistDAWass",dependencies = T) #Installing the first time
2  library(HistDAWass)
```

The main procedures available for the PCA are:

- `WH.1d.PCA` for the analysis of a single distributional variable.
- `WH.MultiplePCA` for analysing more than one distributional variable.

# The Ozone dataset

The dataset contains MatH (matrix of histogram-valued data) object This data set list 78 stations located in the USA recording four variables, without missing data.

```
1  get.MatH.main.info(OzoneFull)
```
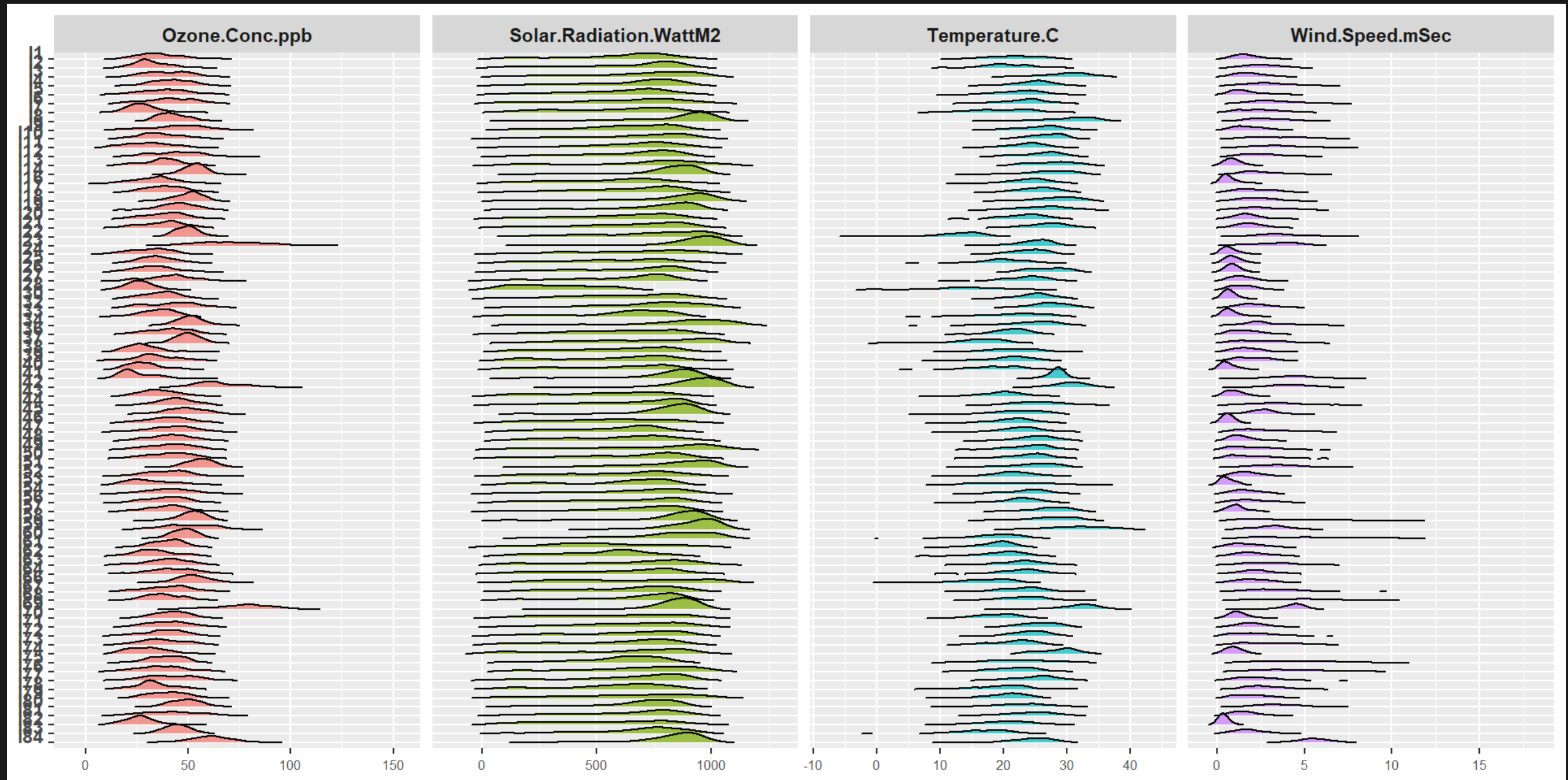
```
$nrows
[1] 78

$ncols
[1] 4

$rownames
 [1] "I1"  "I2"  "I3"  "I4"  "I5"  "I6"  "I7"  "I8"  "I9"  "I10" "I11" "I12"
[13] "I13" "I14" "I16" "I17" "I18" "I19" "I20" "I21" "I22" "I23" "I24" "I25"
[25] "I26" "I27" "I28" "I30" "I32" "I33" "I34" "I36" "I37" "I38" "I39" "I40"
[37] "I41" "I42" "I43" "I44" "I45" "I46" "I47" "I48" "I49" "I50" "I51" "I52"
[49] "I53" "I54" "I56" "I57" "I58" "I59" "I60" "I61" "I62" "I63" "I64" "I66"
[61] "I67" "I68" "I69" "I70" "I71" "I72" "I73" "I74" "I75" "I76" "I77" "I78"
[73] "I79" "I80" "I81" "I82" "I83" "I84"

$varnames
[1] "Ozone.Conc.ppb"        "Temperature.C"         "Solar.Radiation.WattM2"
[4] "Wind.Speed.mSec"
```

# The data

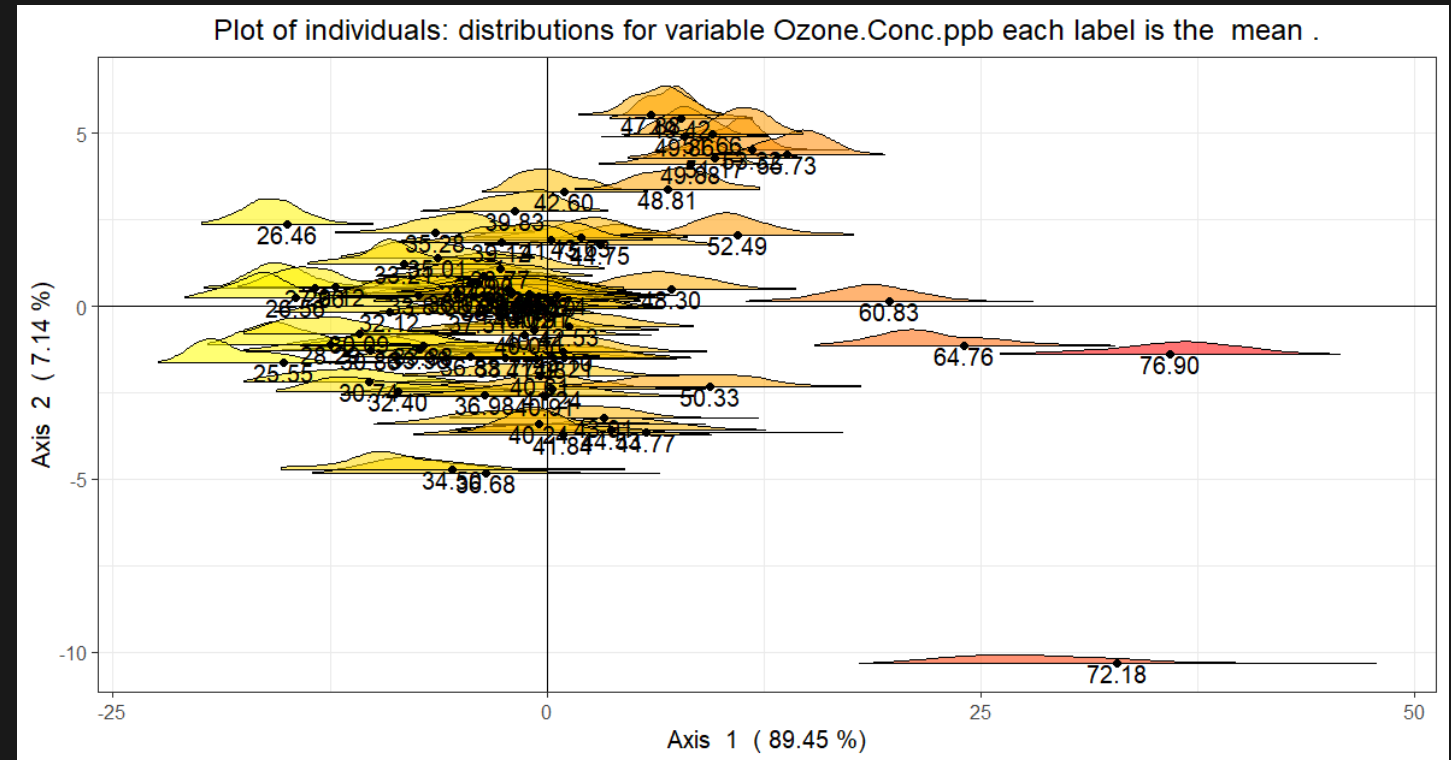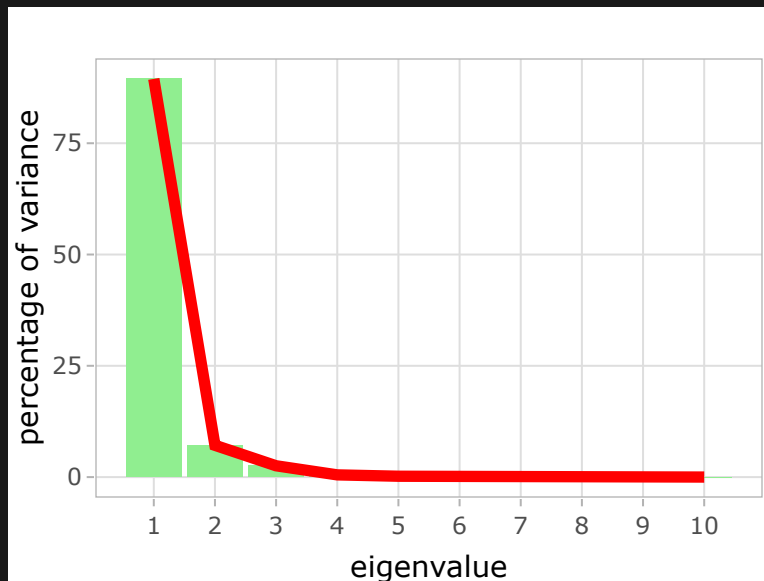# 1d PCA: the analysis of Ozone COncentration ppb

```
1  OZ_1d_PCA<-WH.1d.PCA(OzoneFull,1,quantiles=20)
```
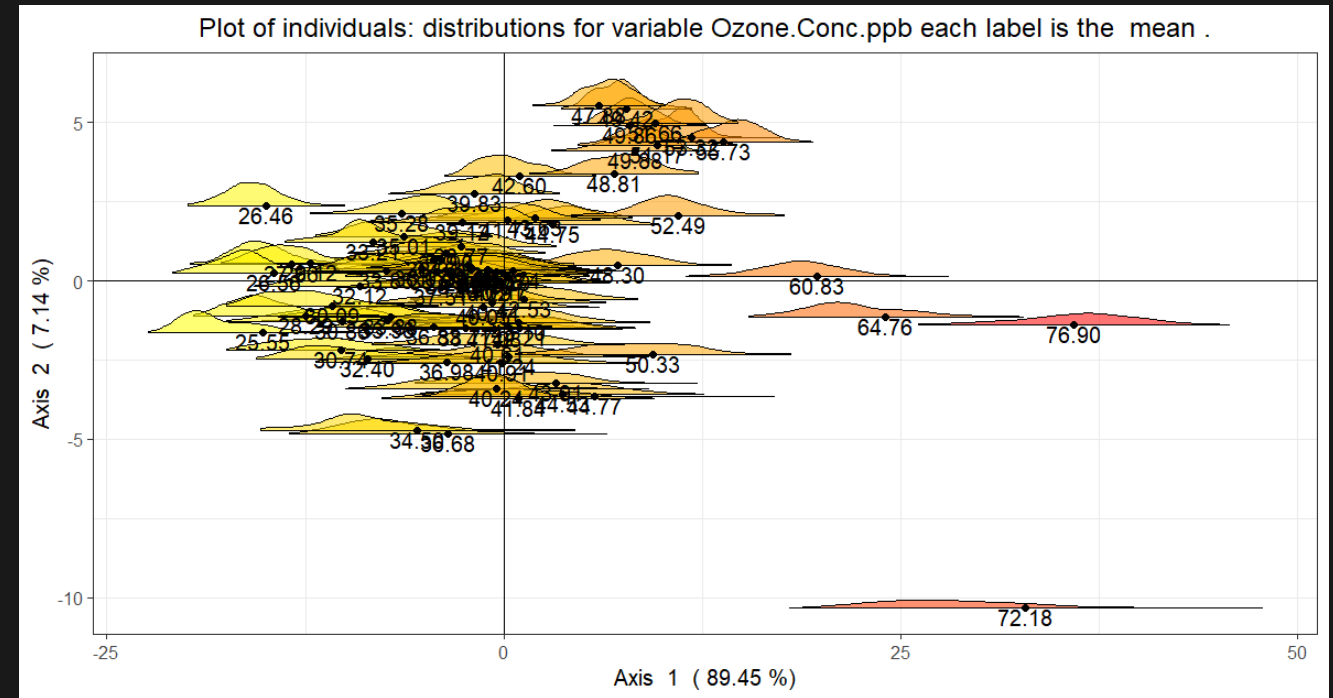
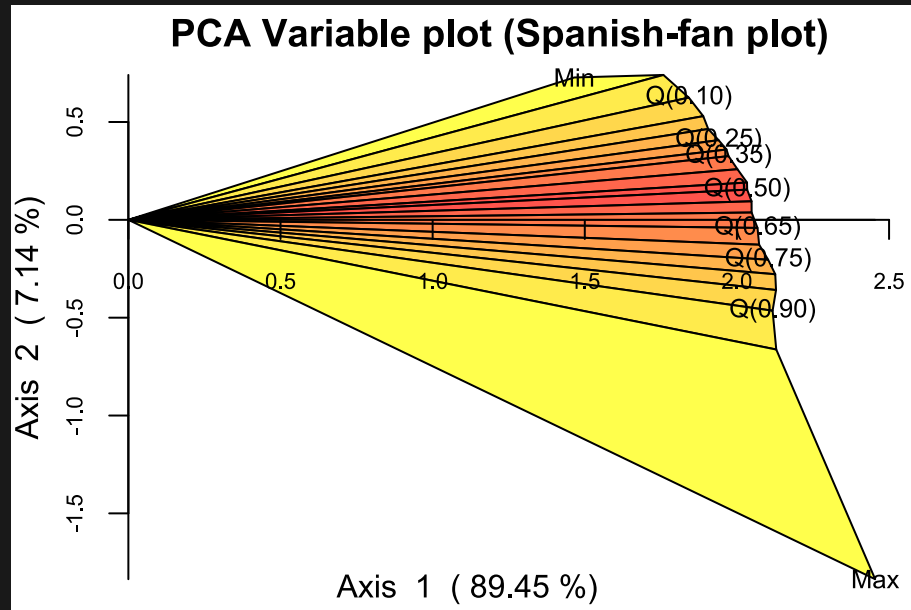We do a PCA on variable --->  Ozone.Conc.ppb

Let's see the scree plot and the first plane





Stations on the first plane

# Interpreting PCs

## The Spanish fan plot for variable correlation





Stations on the first plane

# PCA of all the variables

Now we take into consideration all the four variables and we start performing the multiple PCA using the function `WH.MultiplePCA`

```
1  OZ_PCA<-WH.MultiplePCA(data=OzoneFull,list.of.vars = c(1:4),quantiles = 20)
```
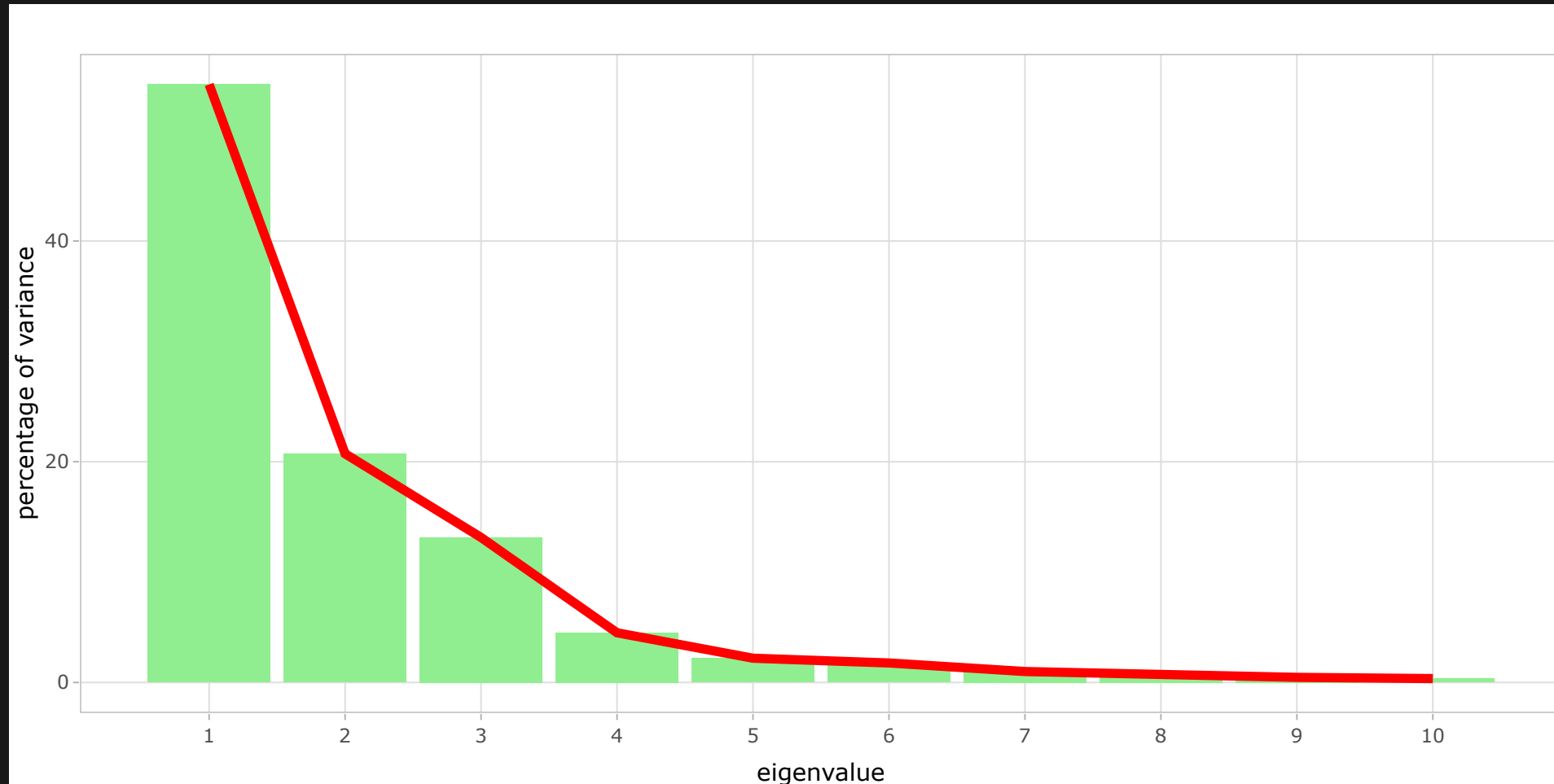
The code executes a Multiple PCA and produces a set of textual and graphical outputs. The code make use of the `FactoMiner` package which is specialized for dimension reduction techniques and the output interpretation.

# Ouput of the procedure

```
**Results of the Multiple Factor Analysis (MFA)**
The analysis was performed on 78 individuals, described by 84 variables
*Results are available in the following objects :

  name                    description
1 "$eig"                  "eigenvalues"
2 "$separate.analyses"    "separate analyses for each group of variables"
3 "$group"                "results for all the groups"
4 "$partial.axes"         "results for the partial axes"
5 "$inertia.ratio"        "inertia ratio"
6 "$ind"                  "results for the individuals"
7 "$quanti.var"           "results for the quantitative variables"
8 "$summary.quanti"       "summary for the quantitative variables"
9 "$global.pca"           "results for the global PCA"
```
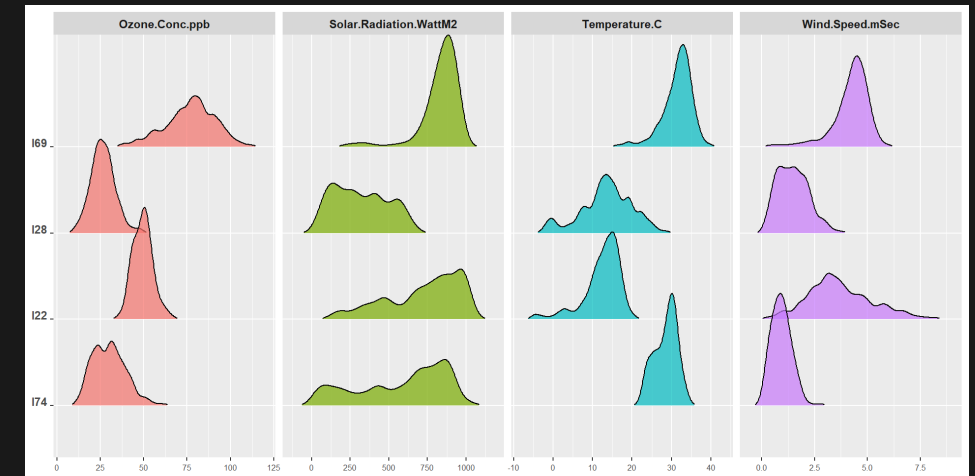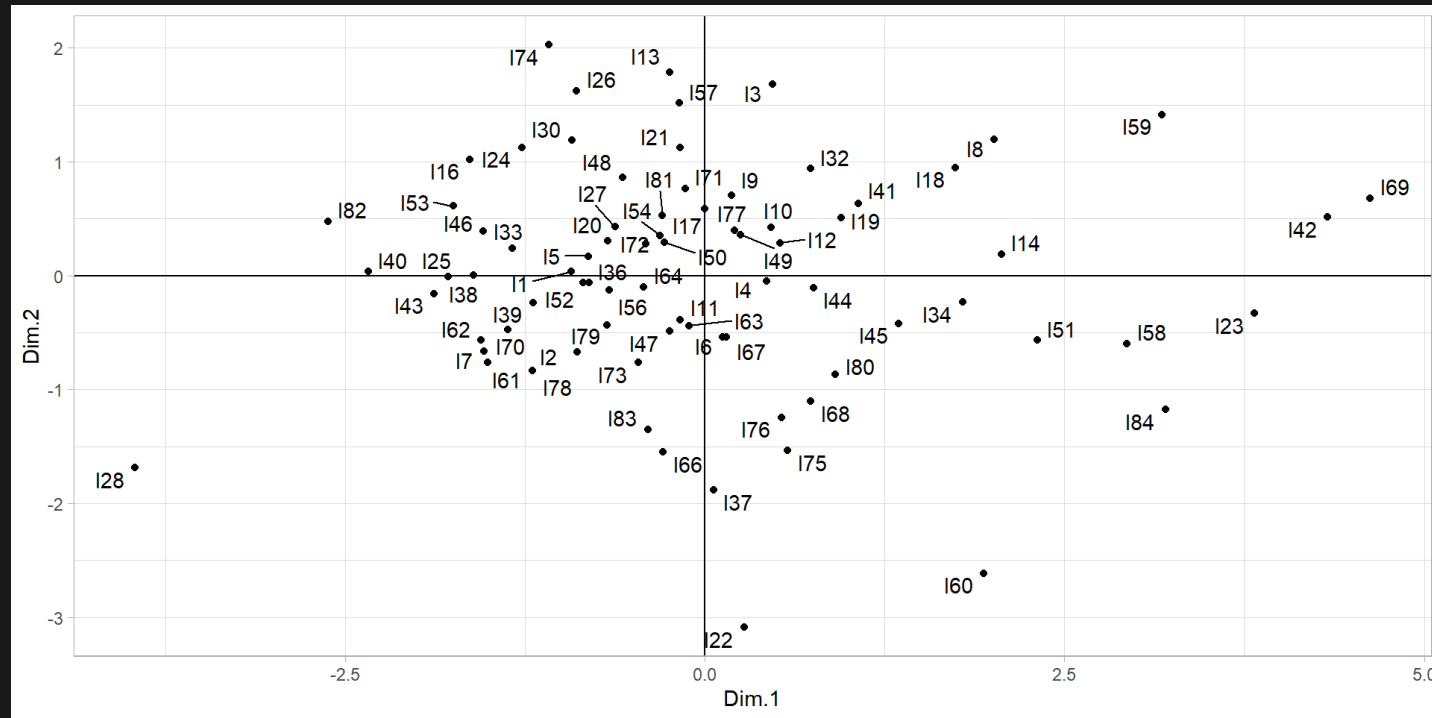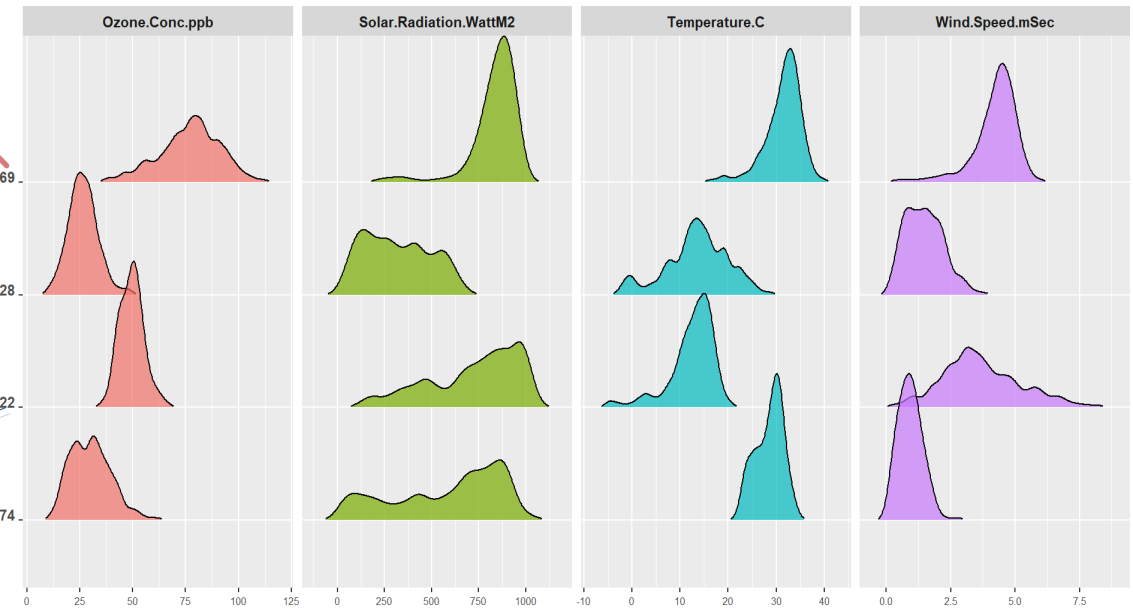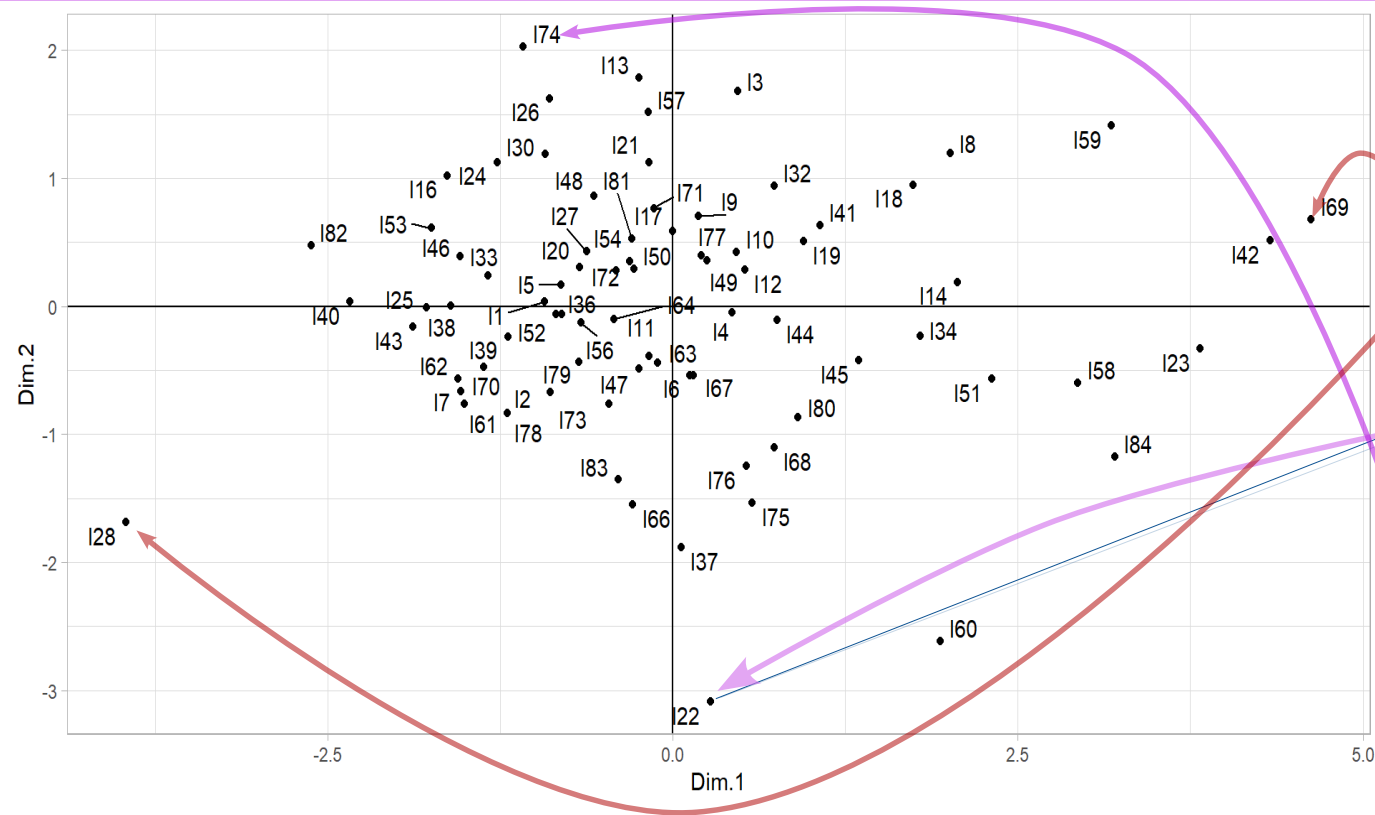
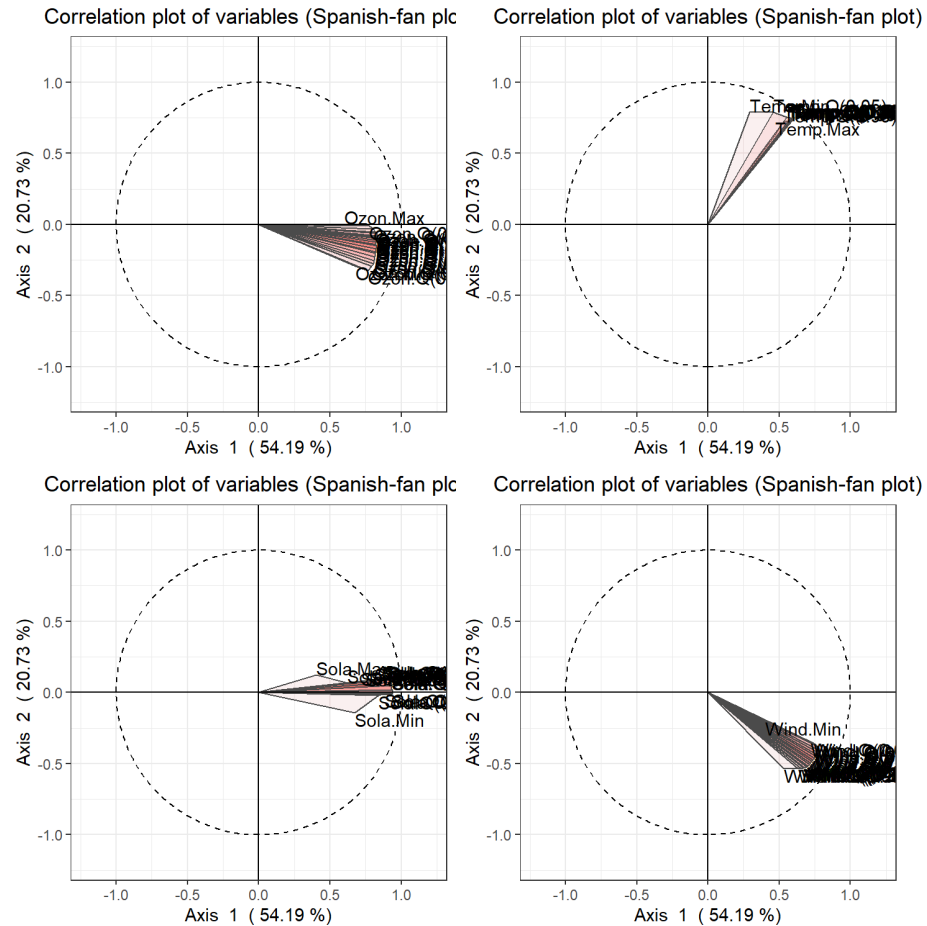# The analysis of the eigenvalues: the scree-plot

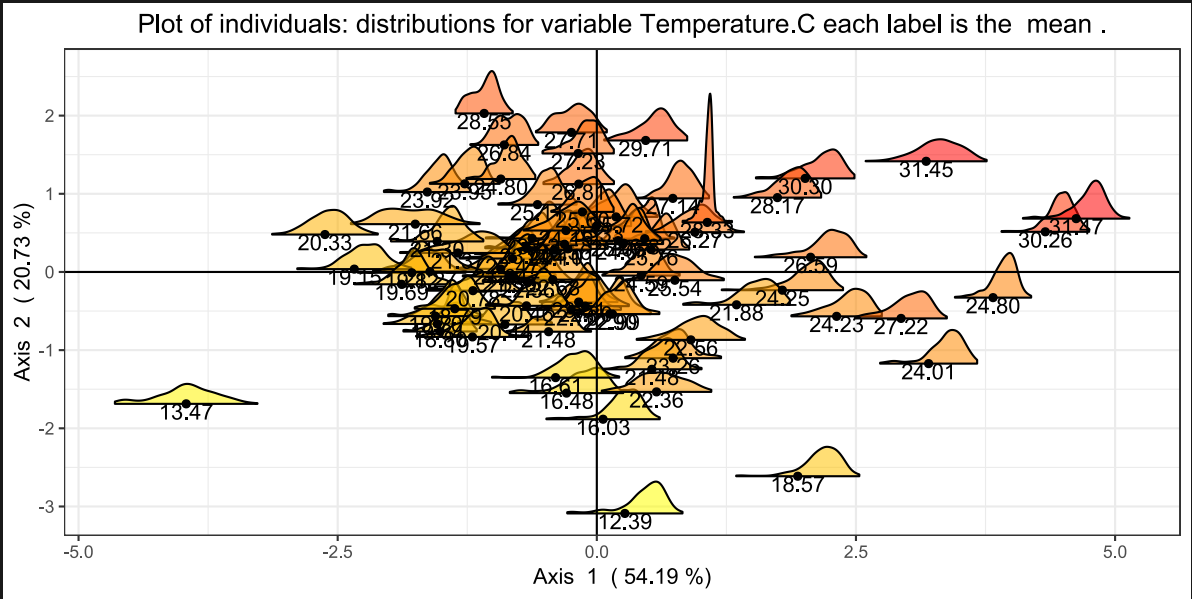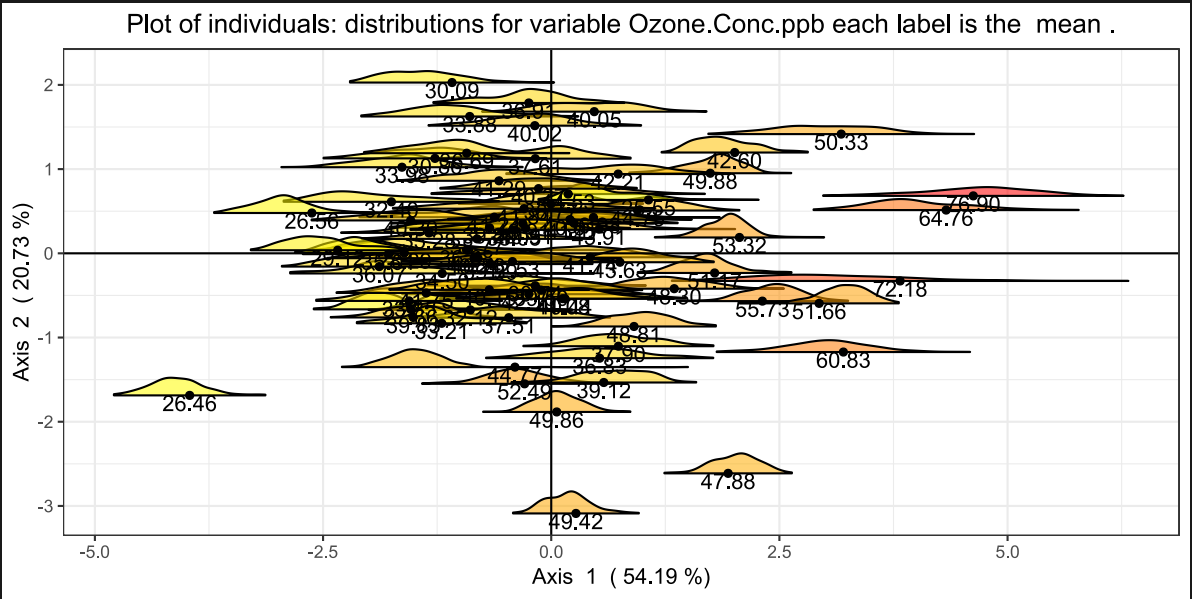# The plot of individuals (the 78 stations) on the first plane

# The plot(s) of variables on the first plane: the spanish fun plots



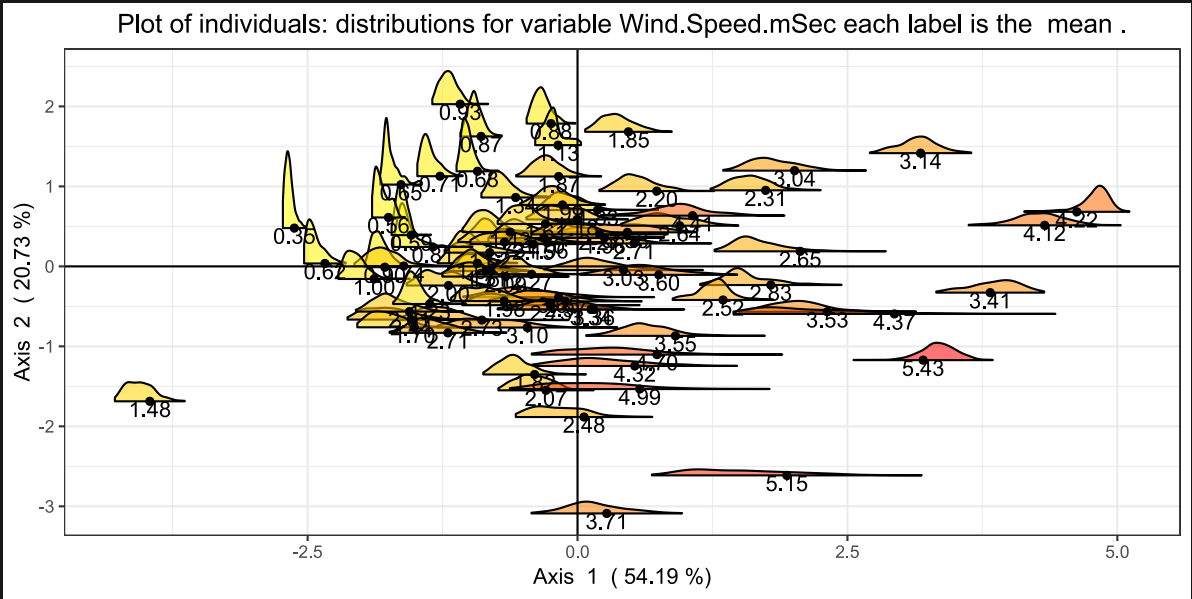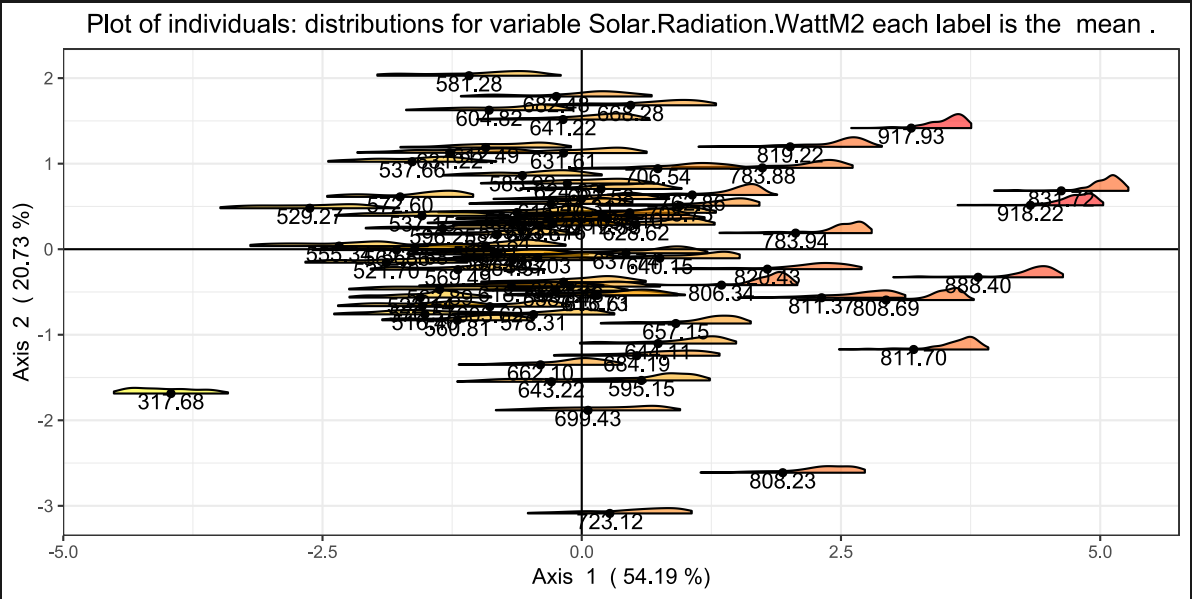Correlation plot of variables (Spanish-fan plot)

**Some comments:**

- The horizontal direction is highly correlated with the Ozone and solar radiation. And both, Ozone and Solar Radiation, are highly and positively correlated.

- The vertical dimension is more related to the Temperature.

- Temperature is rather uncorrelated to the Wind Speed (The two fans are almost at 90 degrees).

- Wind Speed is moderately and positively correlated with the Ozone, but the correlation decrease from the minimum of the wind speed to the maximum.
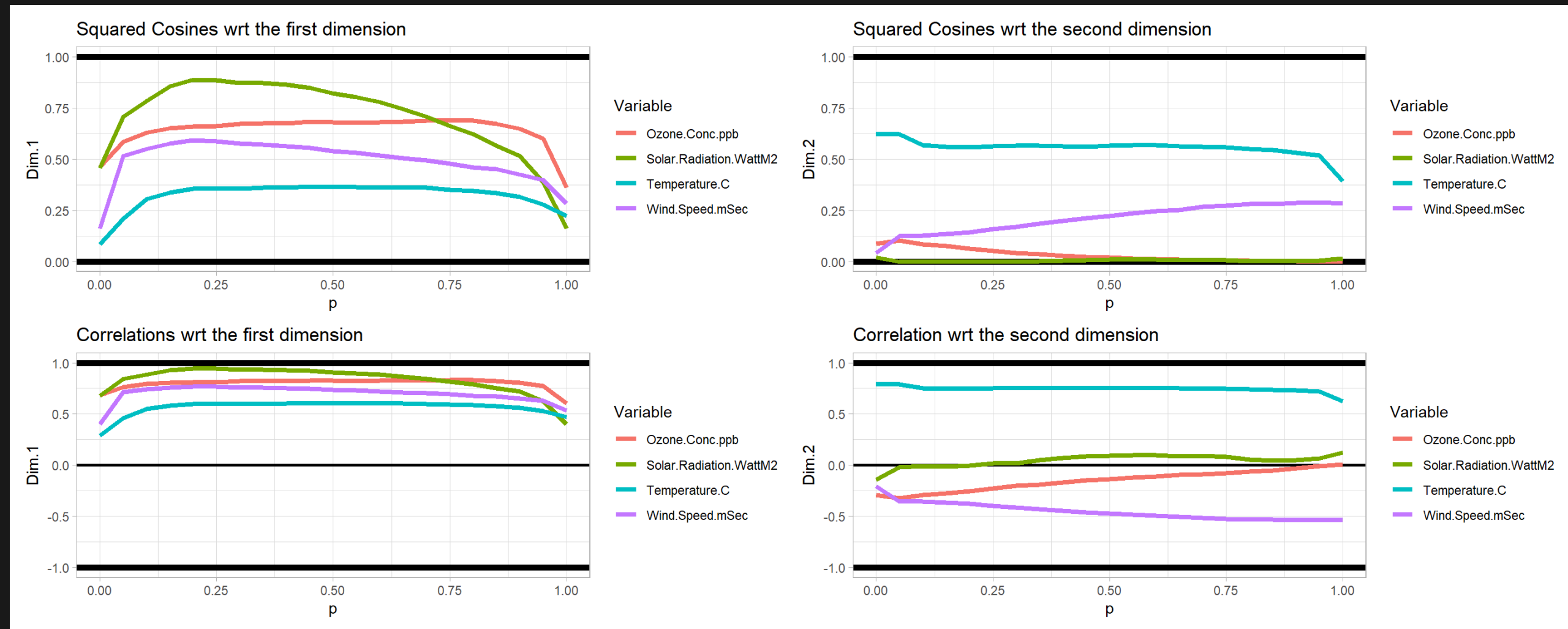
Plot of individuals: distributions for variable Ozone.Conc.ppb each label is the mean .



Plot of individuals: distributions for variable Temperature.C each label is the mean .

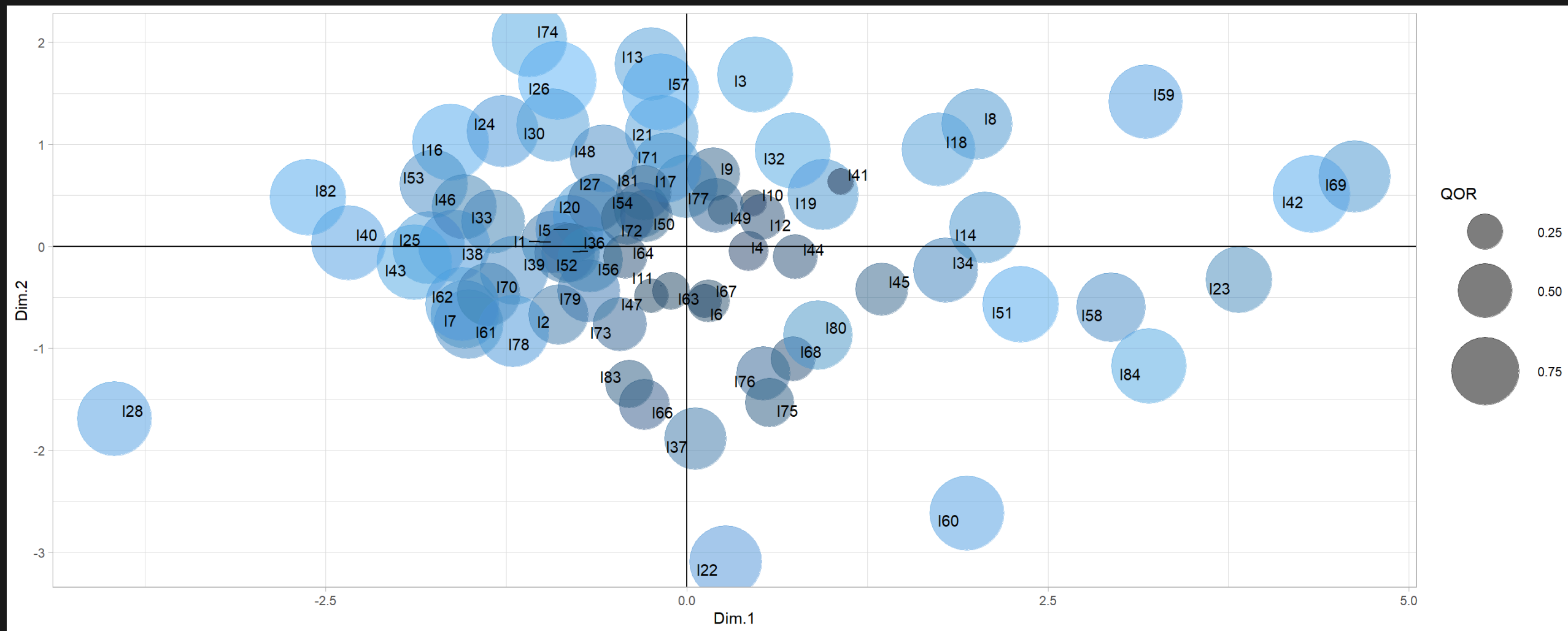# The plot of distributions for each variable 2

# Tools for the interpretation of axes: plots of COS2 and correlations

# The quality of representation of individuals

## The size of the balls is proportional to the quality of representation of points on the plane

# Conclusions

- PCA for distributional data allows to discover more patterns in the data with recpect to the PCA on points;

- If data are intervals, you can consider them as uniform distributions, namely, a histogram with just one bin;

- If distributions have a discrete domain, it is easy to generalize the method.

- If distributions has a nomimal support, you can't use PCA (other methods are available).