

A robust clustering approach to fraud detection

Luis Angel García-Escudero

Dpto. de Estadística e I.O. and IMUVA - Universidad de Valladolid

joint (and on-going...) work with A. Mayo-Isacar, A. Gordaliza, C. Matrán (U. Valladolid) and colleagues from M. Riani, A. Cerioli (U. Parma) and D. Perrotta, F. Torti (JRC-Ispra)

1. CLUSTERING AND ROBUSTNESS

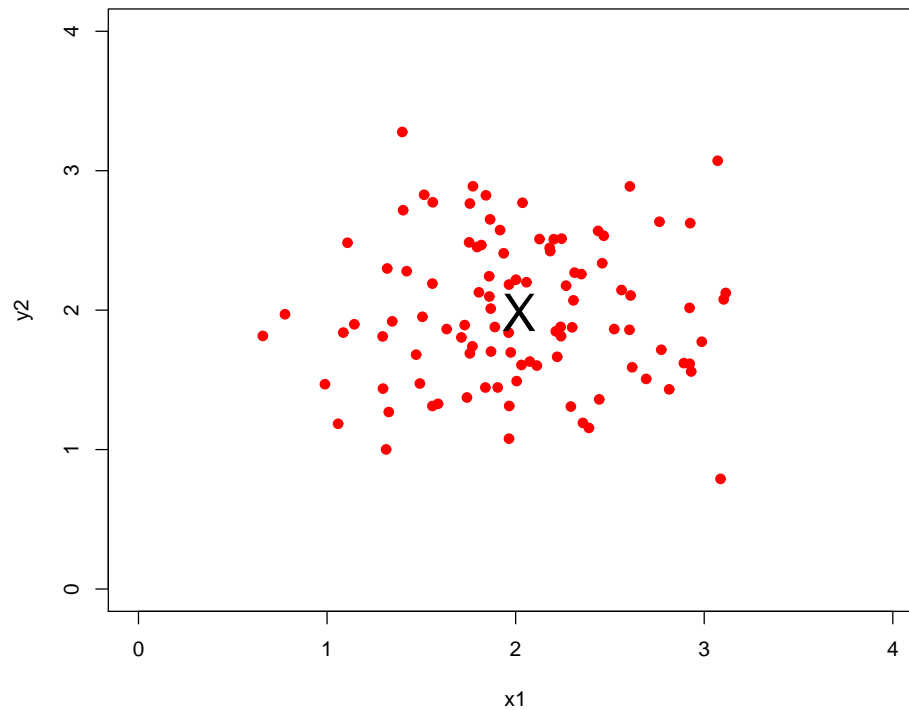
- **Clustering** is the task of **grouping** a set of objects in such a way that objects in the same cluster are more similar to each other than to those in other clusters:



- **Sample mean:**

- ◇ $m = \frac{1}{n} \sum_{i=1}^n x_i$ minimizes $\sum_{i=1}^n \|x_i - m\|^2$

- ◇ m may be seen as the “center” of a data-cloud:



- k clusters $\Rightarrow k$ “data-clouds” $\Rightarrow k$ -means

- **k -means:** Search for
 - ◇ k centers m_1, \dots, m_k
 - ◇ a partition $\{R_1, \dots, R_k\}$ of $\{1, 2, \dots, n\}$

minimizing

$$\sum_{j=1}^k \sum_{i \in R_j} \|x_i - m_j\|^2.$$

- **Cluster j :**

$$R_j = \{i : \|x_i - m_j\| \leq \|x_i - m_l\| \text{ for every } l = 1, \dots, k\}$$

(...assignment to the closest center...)

- **Robustness:** Many statistical procedures are strongly affected by even few outlying observations:

- ◇ The mean is not robust:

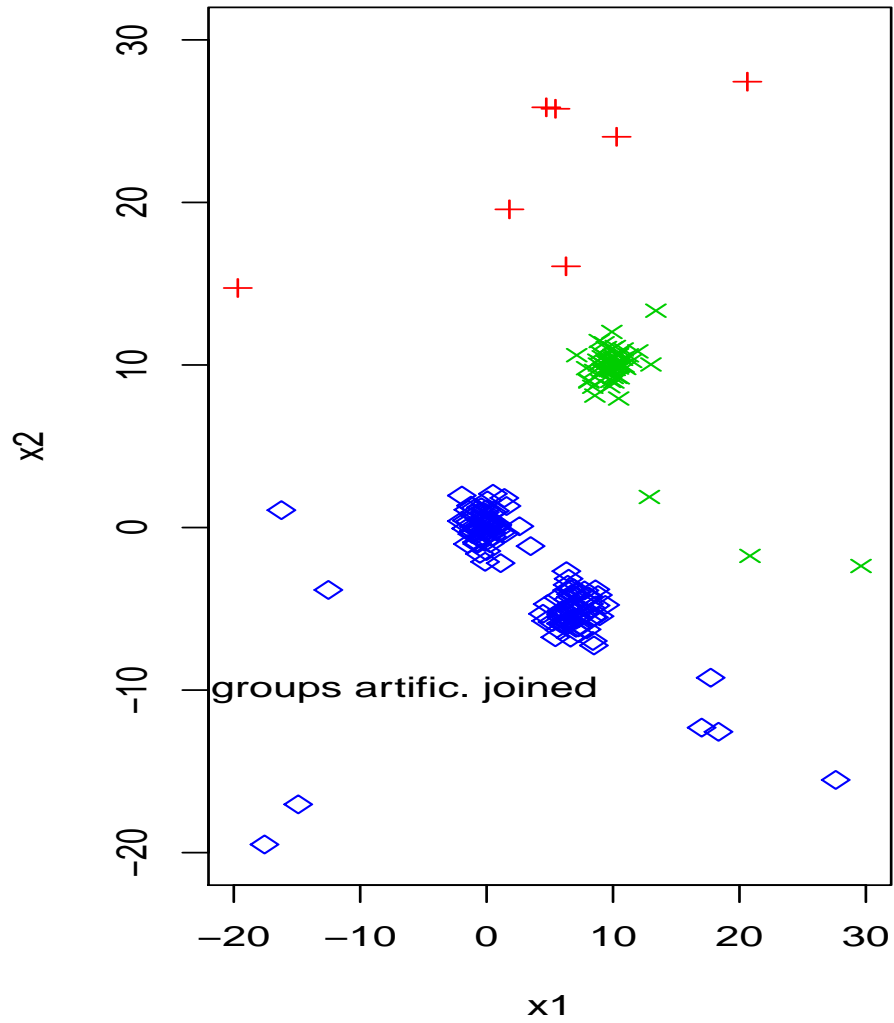
$$\bar{x} = \frac{1.72 + 1.67 + 1.80 + 1.70 + 1.82 + 1.73 + 1.78}{7} = 1.745$$

$$\bar{x} = \frac{1.72 + 1.67 + 1.80 + 1.70 + 182 + 1.73 + 1.78}{7} = 27.485$$

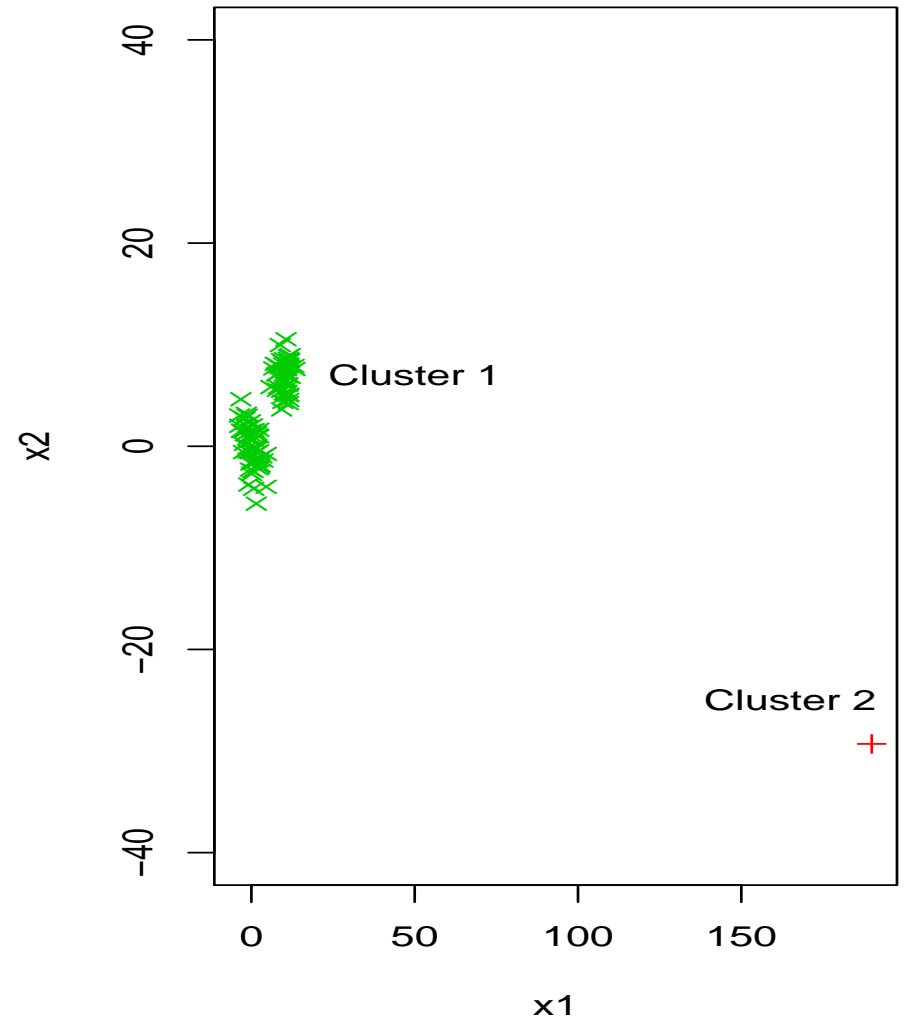
- ◇ *k*-means inherits that **lack of robustness from the mean**

- Lack of robustness of k -means:

(a) 3-means



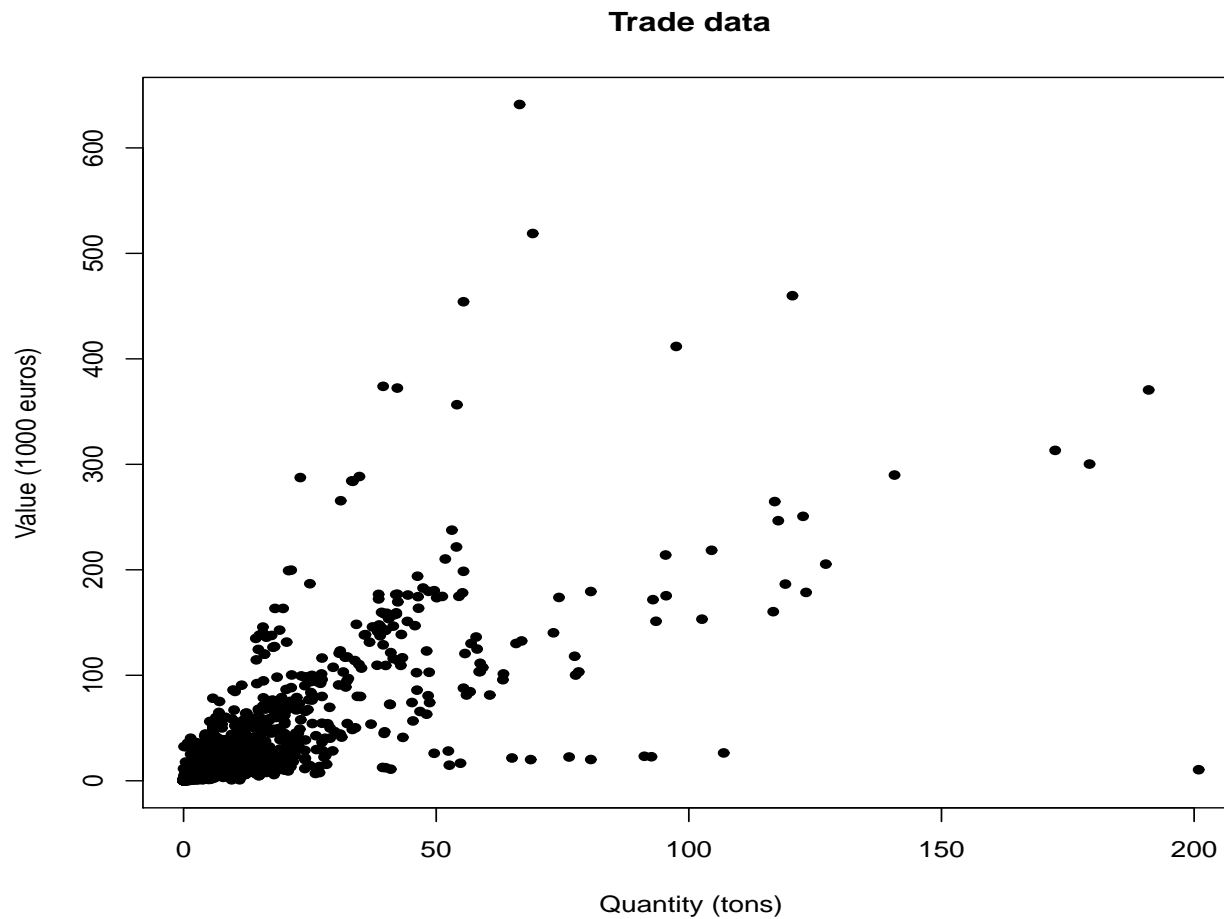
(b) 2-means



- **Outliers** can be seen as “clusters by themselves”
- **So, why not increasing the number of clusters...?**
 - ◇ **But:**
 - Due to (physical, economical,...) reasons we could have an initial idea of k without being aware of the existence of outliers
 - “Radial/background” noise requires large k 's
- **Moreover, the detection of outliers may be the goal itself!!!**

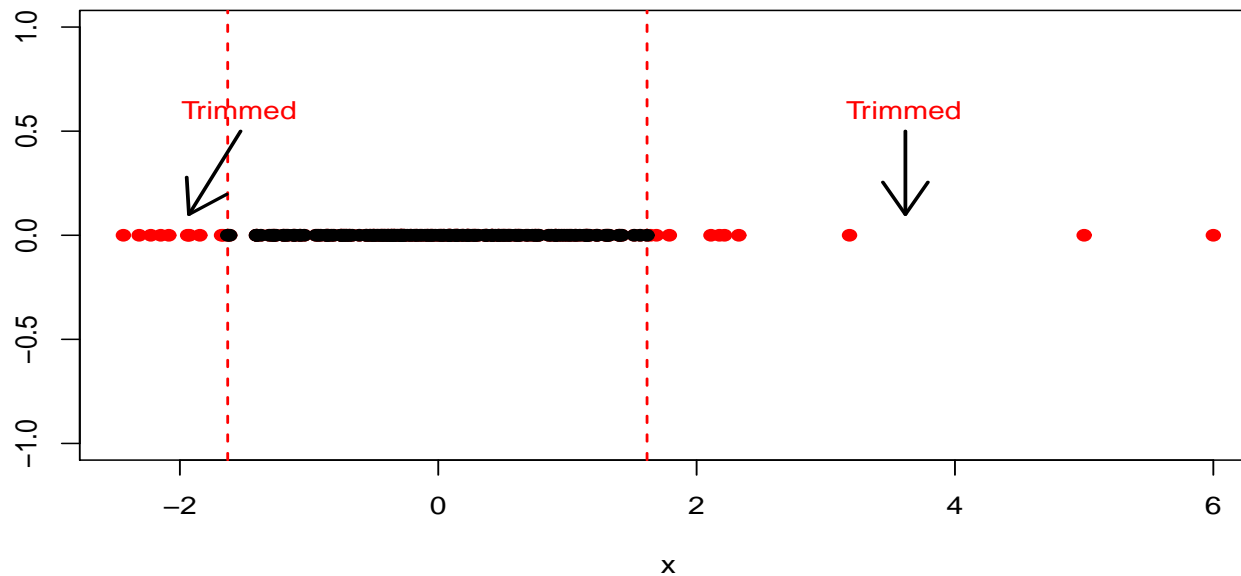
- **Outliers in trade data can be associated to “frauds”:**

- ◇ Heterogeneous sources of data (clustering) + Few outliers (frauds??)



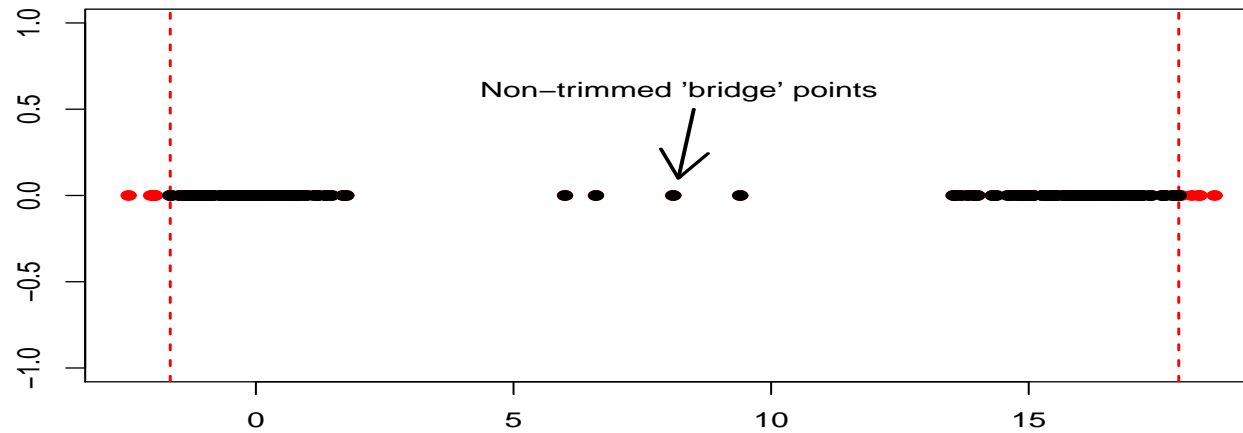
2.- TRIMMED k -MEANS

- **Trimming** is the **oldest** and most widely **used** way to achieve robustness.
- **Trimmed mean:** The proportion $\alpha/2$ smallest and $\alpha/2$ largest observations are discarded before computing the mean:

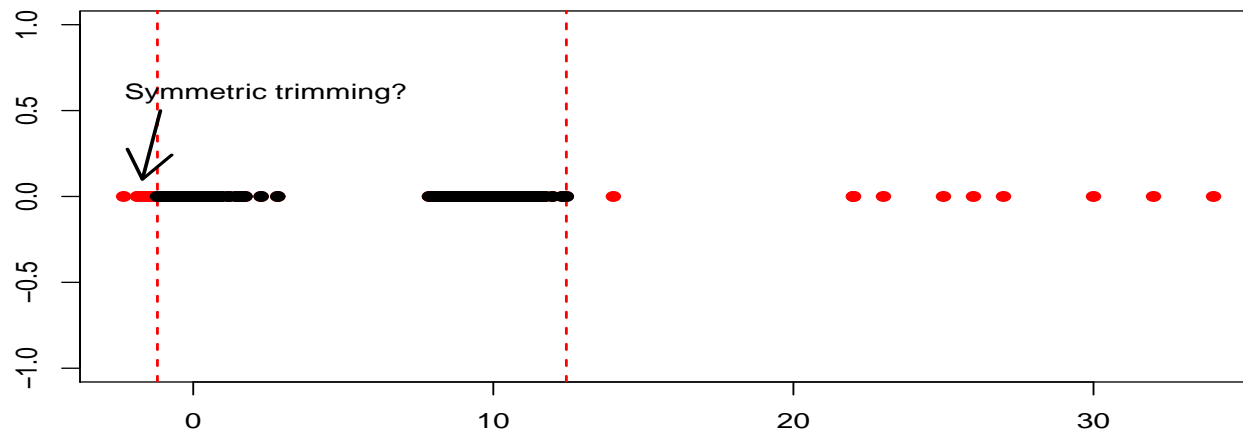


- **But,... how to trim in clustering?**

- ◇ Why not trimming outlying “bridge” points?



- ◇ Why a symmetric trimming?



- ◇ How to trim in *multivariate* clustering problems?

- **Idea: Data itself tell us which are the most outlying observations!!**
 - ◇ Data-driven, adaptive, impartial,... trimming!

- **Trimmed k -means:** we search for

- ◇ k centers m_1, \dots, m_k and

- ◇ a partition $\{R_0, R_1, \dots, R_k\}$ of $\{1, 2, \dots, n\}$ with $\#R_0 = [n\alpha]$

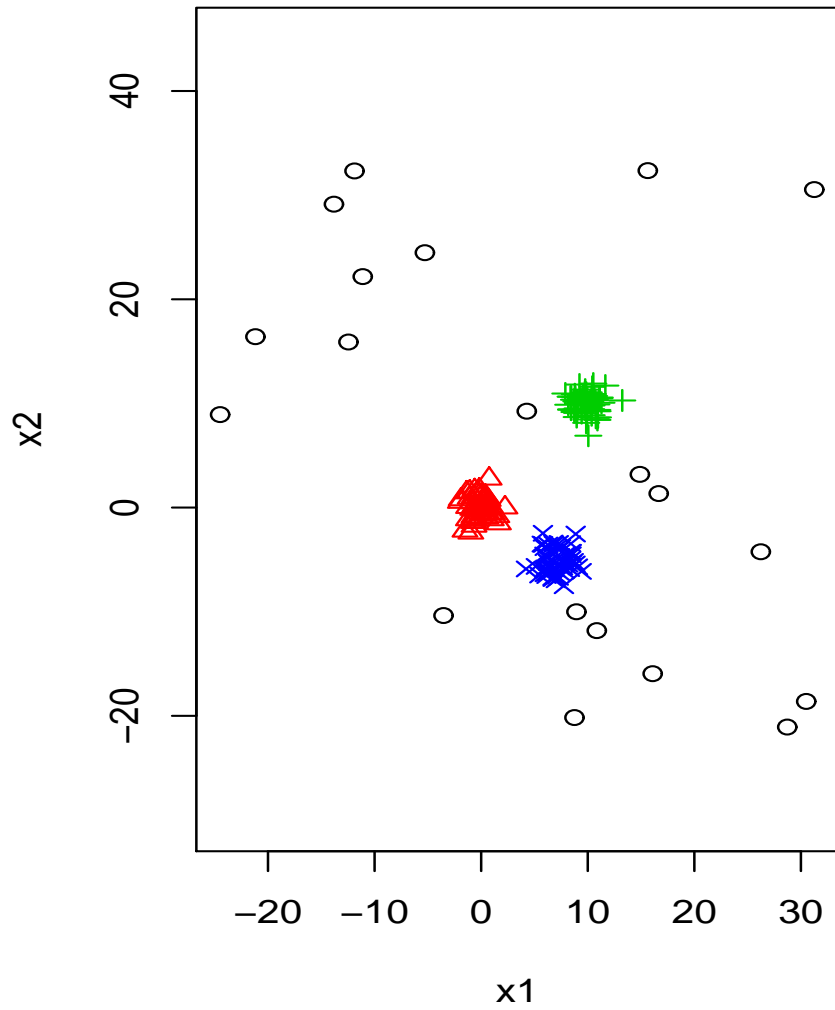
minimizing

$$\sum_{j=1}^k \sum_{i \in R_j} \|x_i - m_j\|^2.$$

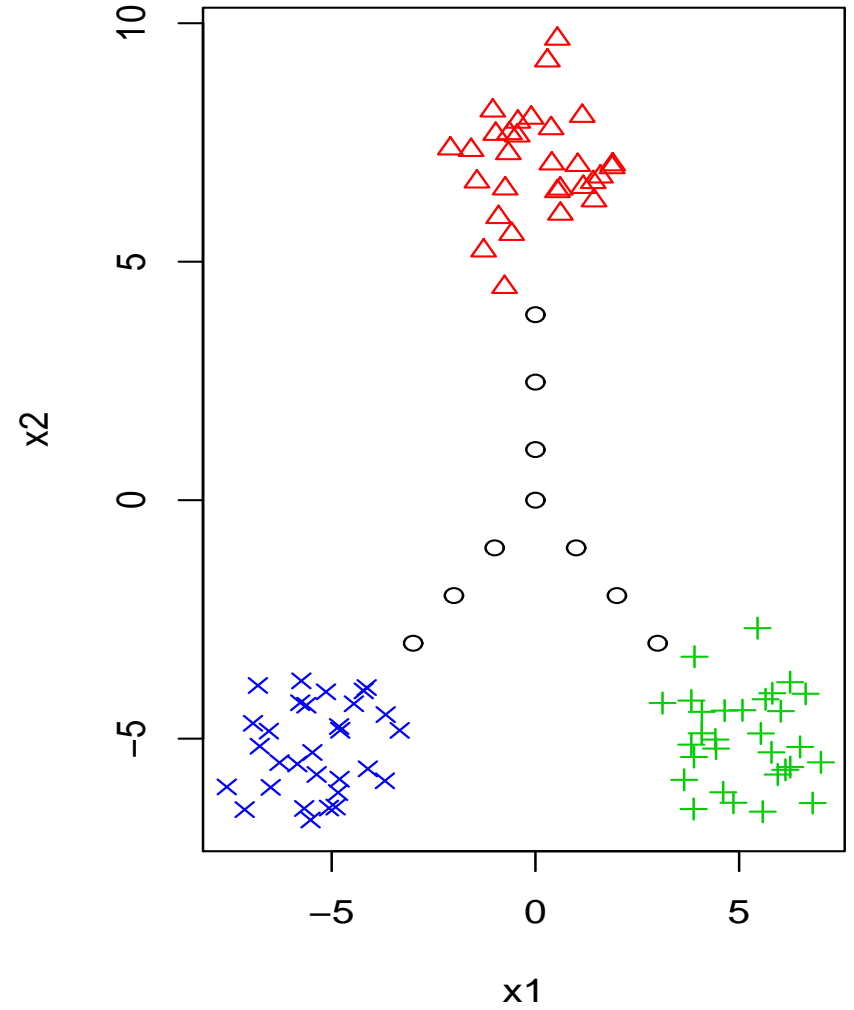
[A fraction α of data is not taken into account \rightsquigarrow Trimmed]

- Black circles: *trimmed points* ($k = 3$ and $\alpha = 0.05$):

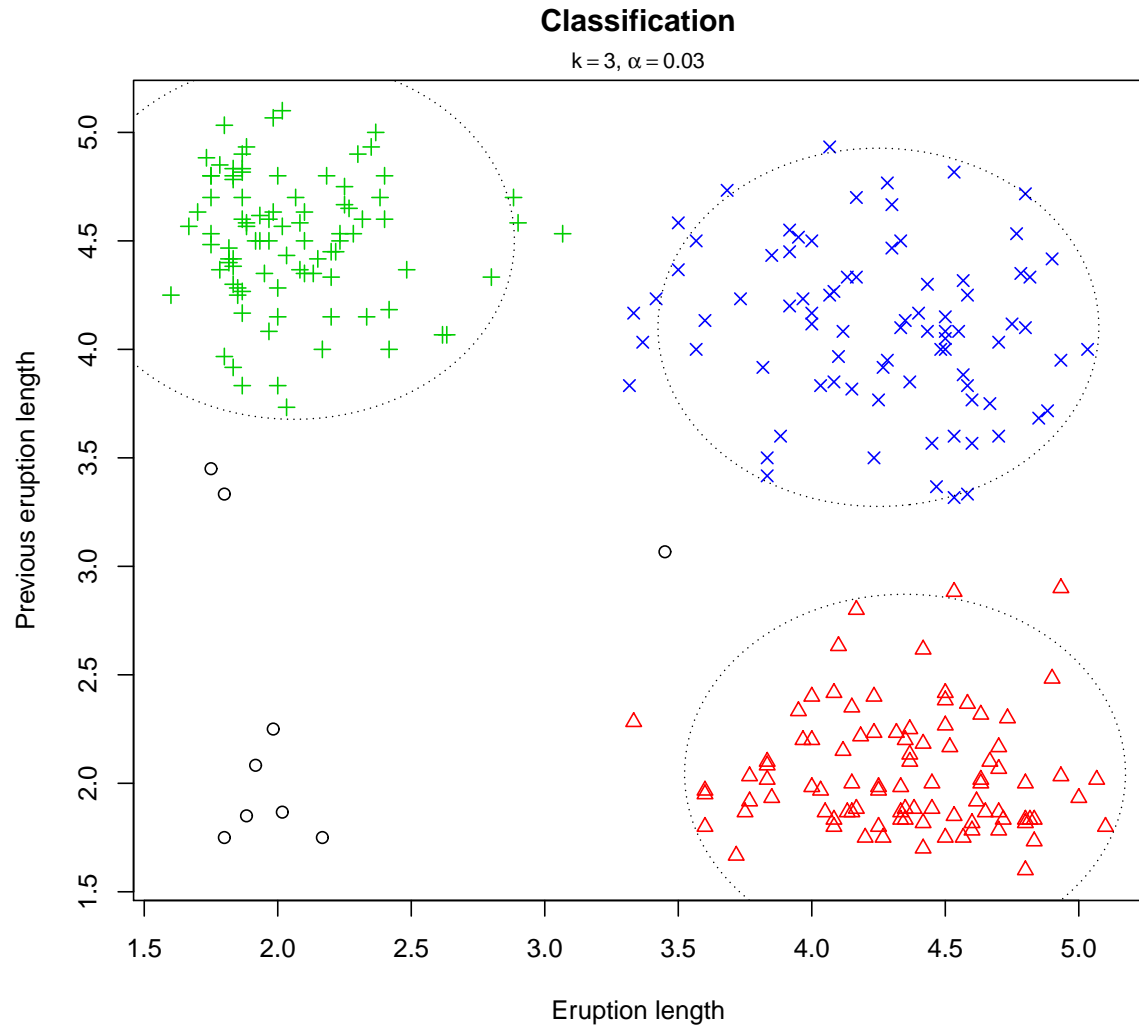
(a)



(b)



- **Old Faithful Geyser data:** $x_1 =$ “Eruption length”, $x_2 =$ “Previous eruption length” and $n = 271$

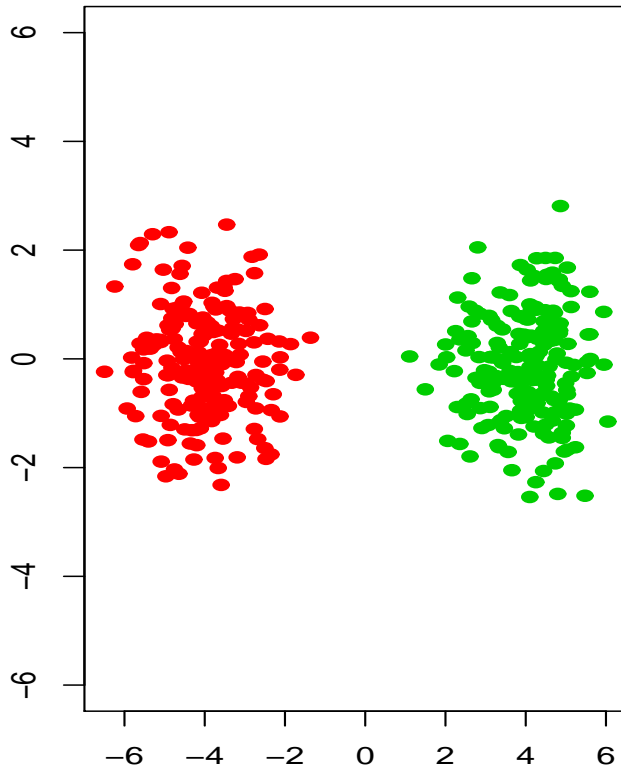


- ◇ $k = 3$ and $\alpha = 0.03$ ($0.03 \cdot 271 \simeq 9$ **trimmed obs.**): 6 rare “short-followed-by-short” eruptions trimmed, 3 bridge points...

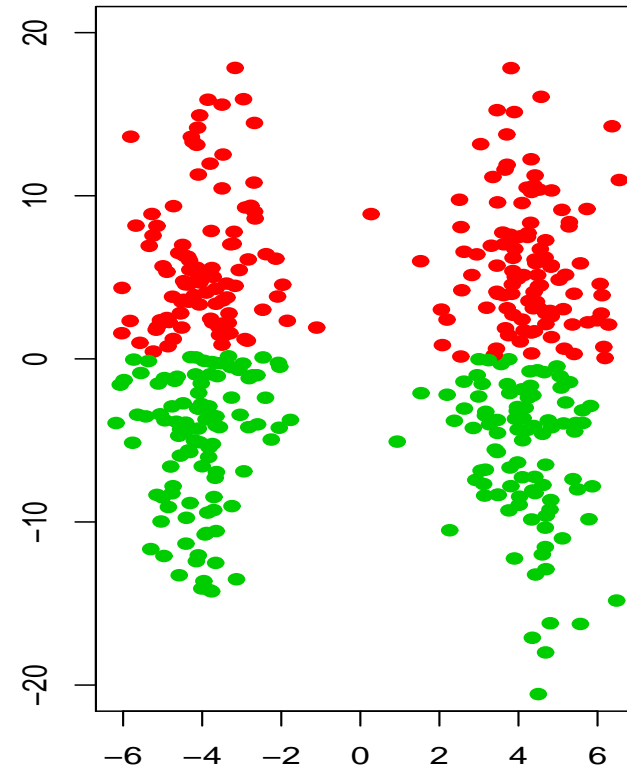
3.- ROBUST MODEL-BASED CLUSTERING

- k -means and trimmed k -means prefer **spherical clusters**:

(a) 2-means (spherical groups)



(b) 2-means (elliptical groups)

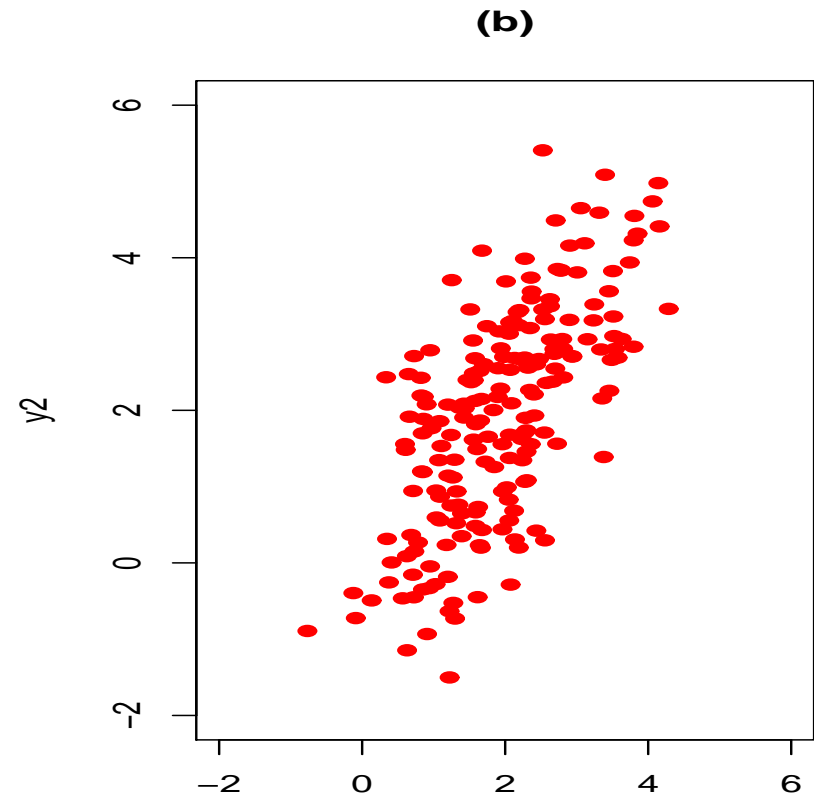
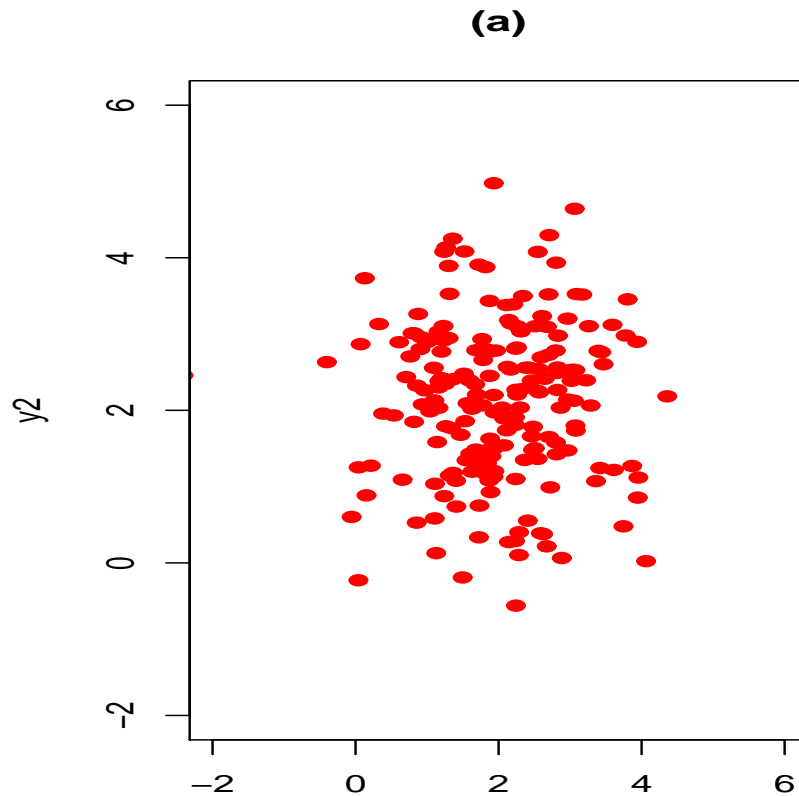


- Elliptically contoured clusters?

- **Multivariate normal** distributions with densities $\phi(\cdot; \mu, \Sigma)$:

- ◇ $\mu = \begin{pmatrix} 2 \\ 2 \end{pmatrix}$ and $\Sigma = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$ [spherical] in (a)

- ◇ $\mu = \begin{pmatrix} 2 \\ 2 \end{pmatrix}$ and $\Sigma = \begin{pmatrix} 2 & 1 \\ 1 & 1 \end{pmatrix}$ [non-spherical] in (b)



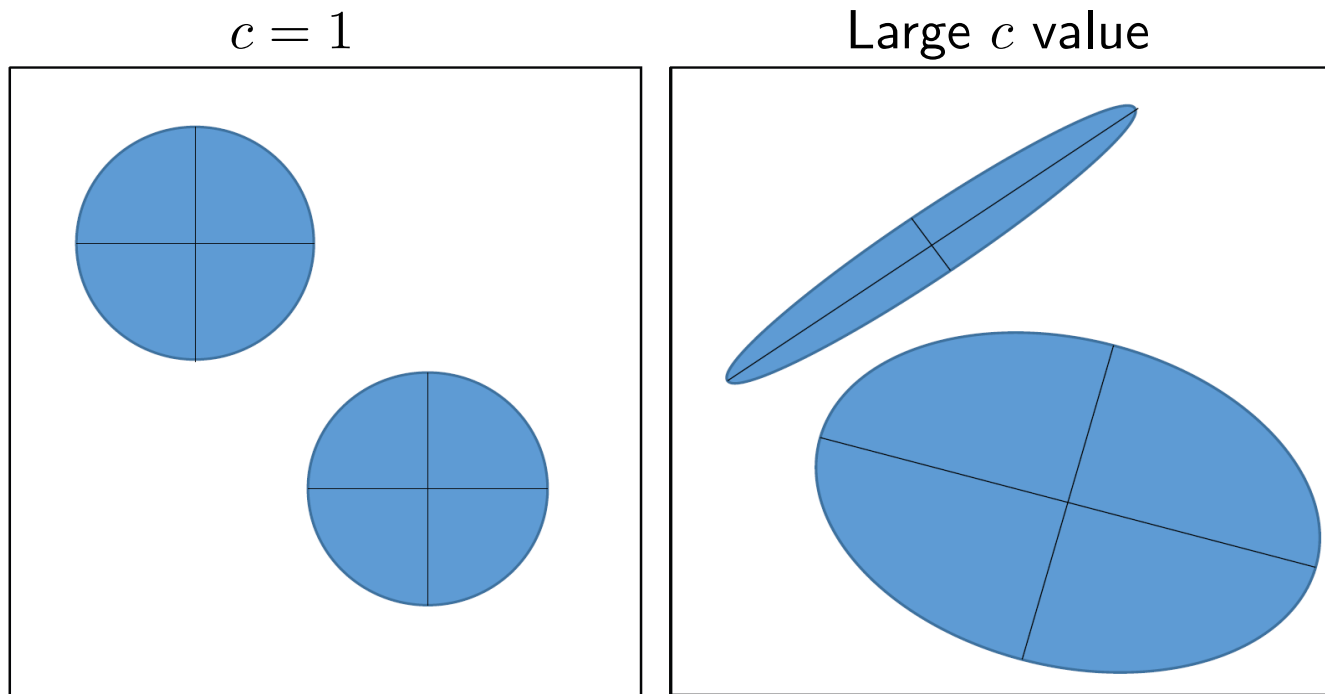
$$\phi(x; \mu, \Sigma) = (2\pi)^{-p/2} |\Sigma|^{-1/2} \exp \left(- (x - \mu)' \Sigma^{-1} (x - \mu) / 2 \right)$$

- **Trimmed likelihoods:** Search for
 - ◇ k centers m_1, \dots, m_k ,
 - ◇ k scatter matrices S_1, \dots, S_k , and,
 - ◇ a partition $\{R_0, R_1, \dots, R_k\}$ of $\{1, 2, \dots, n\}$ with $\#R_0 = [n\alpha]$
 maximizing

$$\sum_{j=1}^k \sum_{x_i \in R_j} \log \phi(x_i; m_j, S_j) \quad (\text{obs. in } R_0 \text{ not taken into account})$$

García-Escudero et al 2008, Neykov et al 2007, Gallegos and Ritter 2005,...

- **Constraints** on the S_j scatter matrices **needed**:
 - ◇ Unbounded target likelihood functions
 - ◇ Avoid detecting (non-interesting) “spurious” clusters
- Control relative axes’ lengths (eigenvalues constraints):



- The **FSDA** Matlab toolbox:

The screenshot shows the documentation for the `tclust` function in the Flexible Statistics and Data Analysis (FSDA) toolbox. The page includes a search bar, a navigation menu, and the following content:

tclust
tclust computes trimmed clustering with scatter restrictions

Syntax

```
out = tclust(Y,k,alpha,restfactor)
out = tclust(Y,k,alpha,restfactor,Name,Value)
[out , varargout]=tclust(___)
```

Description

tclust partitions the points in the n-by-v data matrix Y into k clusters. This partition minimizes the trimmed sum, over all clusters, of the within-cluster sums of point-to-cluster-centroid distances. Rows of Y correspond to points, columns correspond to variables. tclust returns inside structure out an n-by-1 vector idx containing the cluster indices of each point. By default, tclust uses (squared), possibly constrained, Mahalanobis distances.

`out = tclust(Y, k, alpha, restfactor)` tclust of geysers data using k=3, alpha=0.1 and restfactor=10000.

`out = tclust(Y, k, alpha, restfactor, Name, Value)` Use of 'plots' option as a struct, to produce more complex plots.

`[out , varargout] = tclust(___)` tclust of geysers with varargout.

- The R package **tclust** at CRAN repository:

The screenshot shows the RGui interface with the R Console and the documentation for the `tclust` package. The R Console shows the following commands and output:

```
R version 2.10.0 (2009
Copyright (C) 2009 The
ISBN 3-900051-07-0

R es un software libre
Usted puede redistribu
Escriba 'license()' o

R es un proyecto colab
Escriba 'contributors('
'citation()' para sabe
Escriba 'demo()' para
o 'help.start()' para
Escriba 'q()' para sal

[Previously saved work
> library(tclust)
Loading Required packa
Mensajes de aviso perd
1: package 'tclust' wa
2: package 'mvtnorm' w
> help(package=tclust)
>
```

The documentation window shows the following information:

Information on package 'tclust'

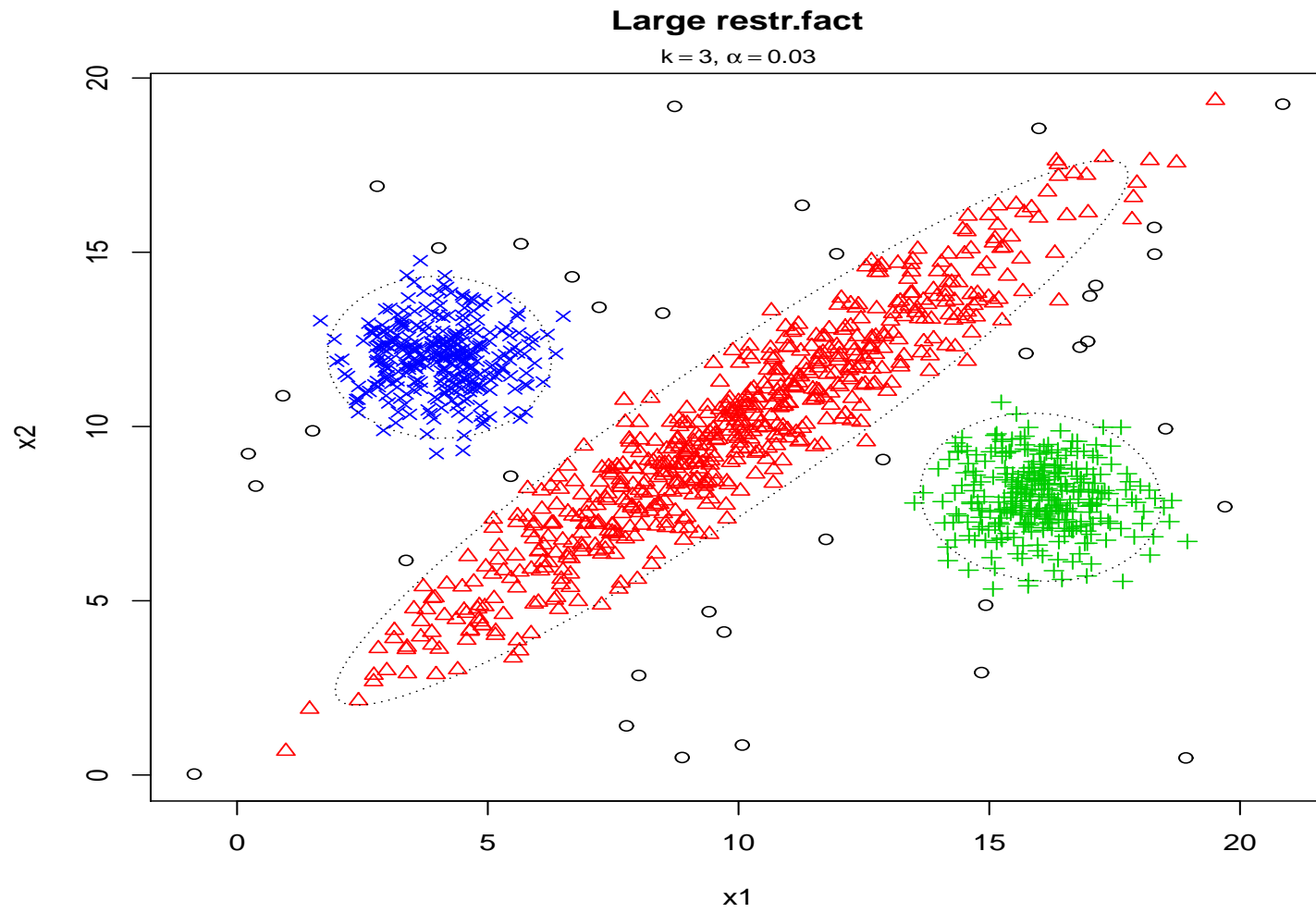
Description:

```
Package:          tclust
Type:             Package
Title:            Robust Trimmed Clustering
Version:          1.0-3
Date:             2009-05-13
Author:           Agustín Mayo Iscar, Luis Ángel García Escudero,
                  Heinrich Frits
Maintainer:       Heinrich Frits <Heinrich.Frits@hotmail.com>
Description:     Robust Trimmed clustering
License:          GPL-3
Depends:          mvtnorm
Packaged:         2010-05-01 09:49:04 UTC: Max
Repository:       CRAN
Date/Publication: 2010-05-01 16:00:22
Built:            R 2.10.1; 1386-pc-mingw32; 2010-05-02 15:21:23 UTC;
                  windows

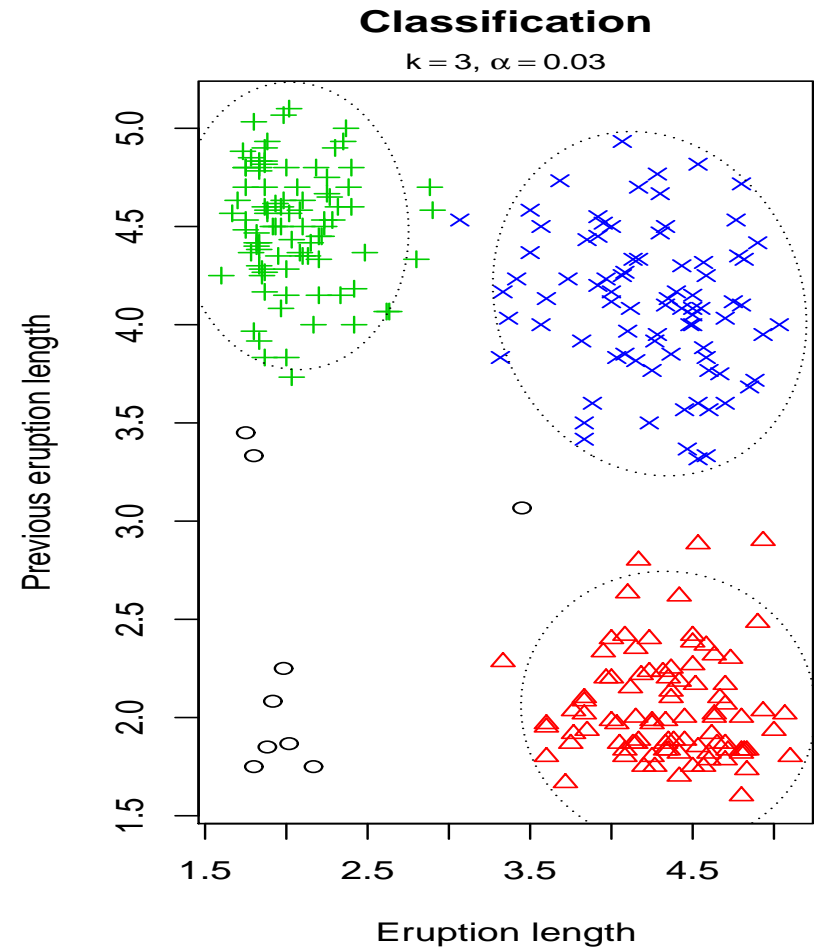
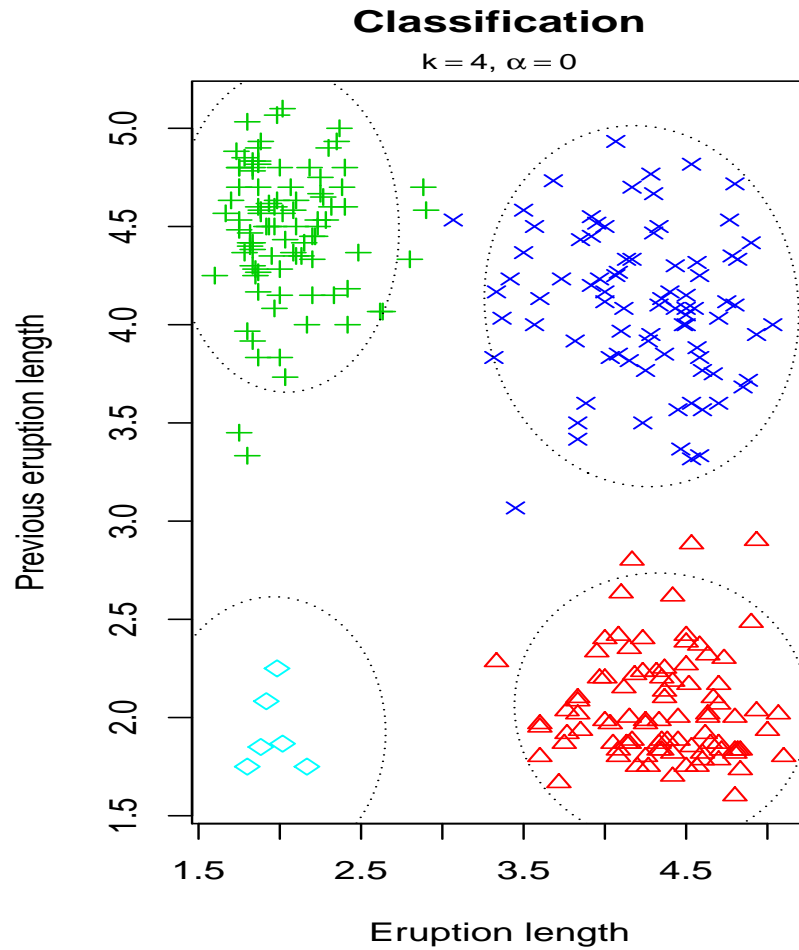
Index:
DiscrFact:       Discriminant Factor Analysis for tclust Objects
```

- The R package `tclust`:
 - > `library(tclust)`
- `tkmeans(data, k, alpha)`
 - ◇ `k` = “number of groups”
 - ◇ `alpha` = “trimming proportion”
- `tclust(data, k, alpha, restr.fact, ...)`
 - ◇ `restr.fact` = “Strength of the constraints”

- `tclust(X,k=3,alpha=0.03,restr.fact=50)`

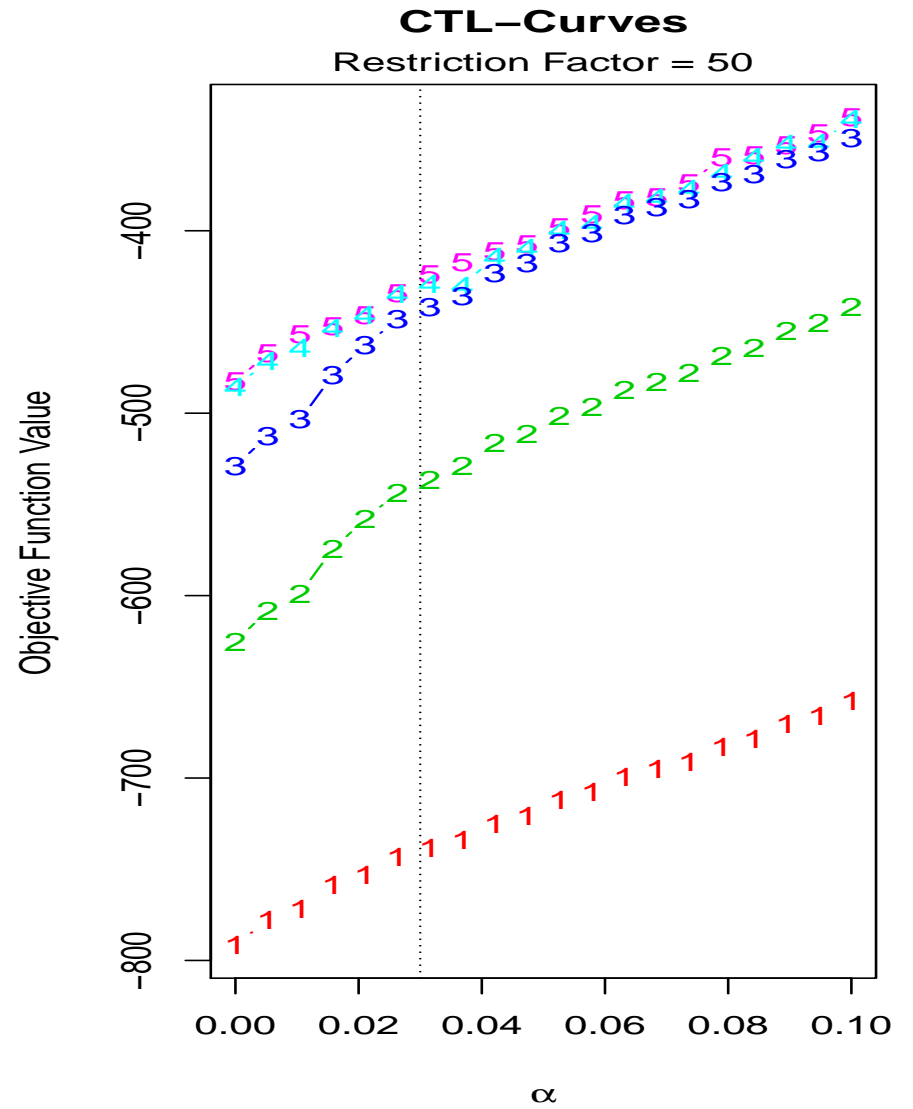
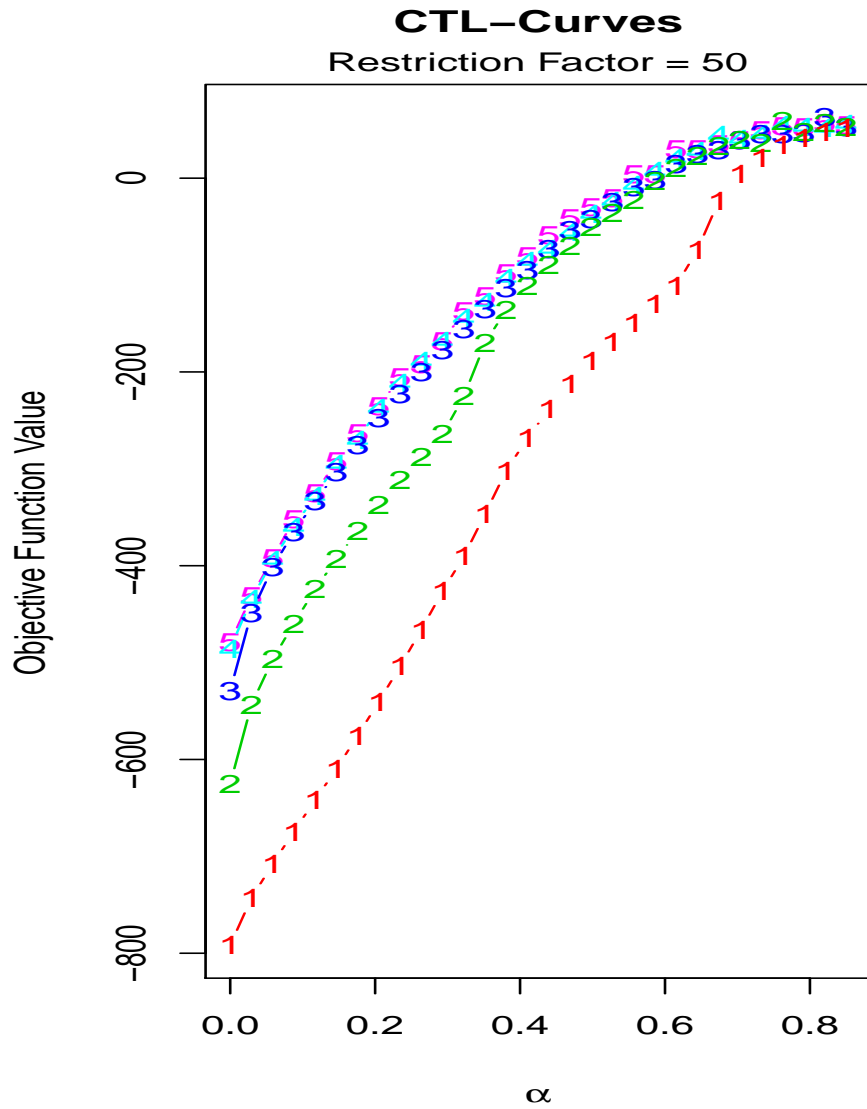


- Old Faithful Geyser data again:



- Why $k = 3$ and $\alpha = 0.03$ was a sensible solution?

- Applying `ctlcurves` to the Old Faithful Geyser data:



4.- ROBUST CLUSTERING AROUND LINEAR SUBSPACES

- **Robust linear grouping:** Higher p dimensions, but assuming that our data “live” in k low-dimensional (affine) subspaces...

◇ We search for

· k linear subspaces h_1, \dots, h_k in \mathbb{R}^p

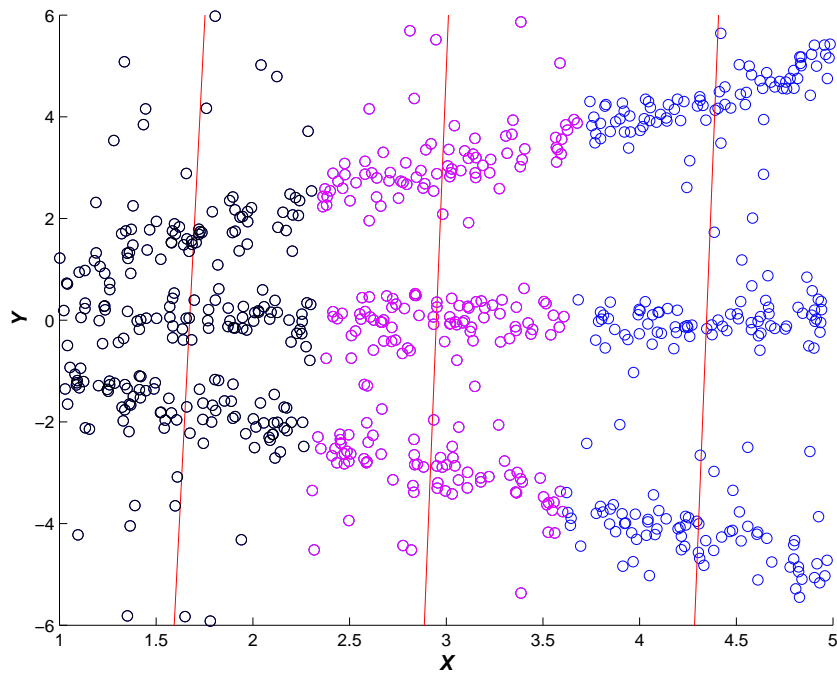
· a partition $\{R_0, R_1, \dots, R_k\}$ of $\{1, 2, \dots, n\}$ with $\#R_0 = [n\alpha]$

minimizing

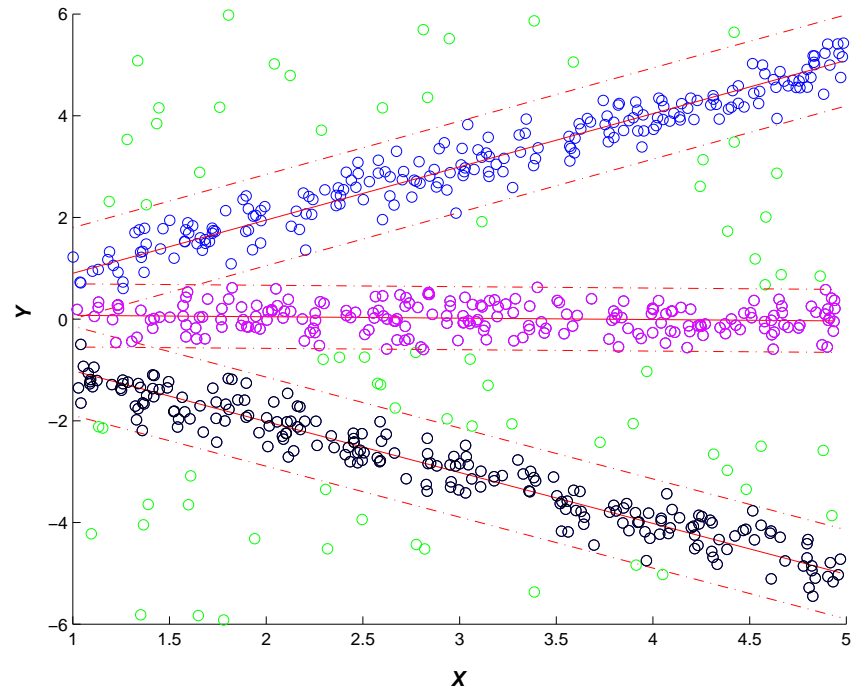
$$\sum_{j=1}^k \sum_{i \in R_j} \|x_i - \text{Pr}_{h_j}(x_i)\|^2.$$

◇ $\text{Pr}_h(\cdot)$ denotes the “orthogonal” projection onto the linear subspace h

- **Example:** Three linear structures in presence of noise:



(a) $\alpha = 0$



(b) $\alpha = 0.1$ (\circ = "Trimmed")

Trimmed "mixtures of regressions" can also be applied...

- $k = 1$ case \Rightarrow **Robust “Principal Components Analysis (PCA)”**:

- ◇ PCA provides a q -dimensional ($q \ll p$) representation of data by

$$\min_{\mathbf{B}_q, \mathbf{A}_q, m} \sum_{i=1}^n \|\mathbf{x}_i - \hat{\mathbf{x}}_i\|^2 \text{ for}$$

$$\hat{\mathbf{x}}_i = \text{Pr}_h(\mathbf{x}_i) = \hat{\mathbf{x}}_i(\mathbf{B}_q, \mathbf{A}_q, m) = m + \mathbf{B}_q \mathbf{a}_i$$

- $\mathbf{A}_q = \begin{pmatrix} -\mathbf{a}_1- \\ \dots \\ -\mathbf{a}_i- \\ \dots \\ -\mathbf{a}_n- \end{pmatrix}$ is the **scores** matrix ($n \times q$)

- $\mathbf{B}_q = \begin{pmatrix} -\mathbf{b}_1- \\ \dots \\ -\mathbf{b}_j- \\ \dots \\ -\mathbf{b}_p- \end{pmatrix}$ is a matrix ($p \times q$) whose columns generate a q -dimensional **approximating subspace** h

- **Principal Components Analysis is highly non-robust!!!**
- **Least Trimmed Squares PCA** (Maronna 2005): Minimize

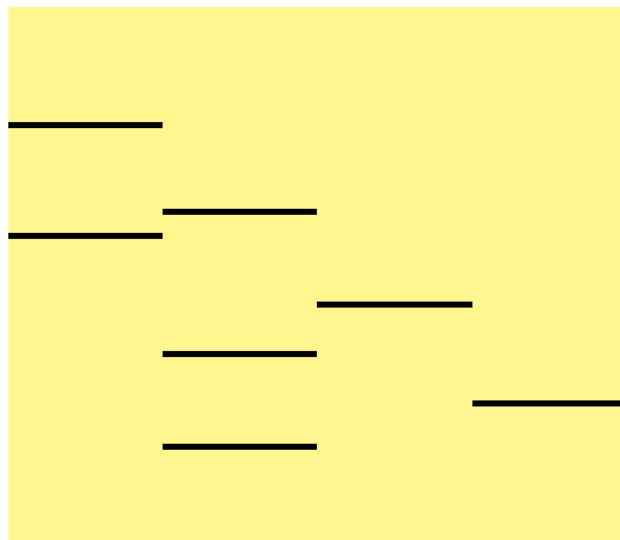
$$\sum_{i=1}^n w_i \|\mathbf{x}_i - \hat{\mathbf{x}}_i\|^2 = \sum_{i=1}^n w_i \|\mathbf{x}_i - \hat{\mathbf{x}}_i(\mathbf{B}_q, \mathbf{A}_q, \mathbf{m})\|^2,$$

with $\{w_i\}_{i=1}^n$ being “0-1 **weights**” such that

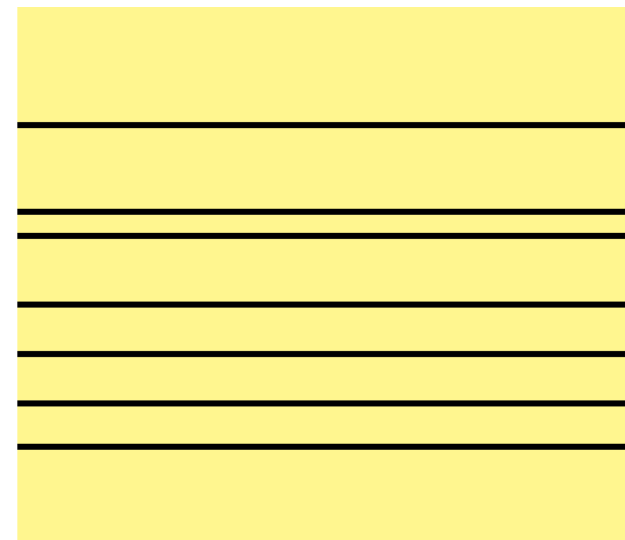
$$\sum_{i=1}^n w_i = [n(1 - \alpha)]$$

◇ **Weights:** $w_i = \begin{cases} 1 & \text{If } \mathbf{x}_i \text{ is not trimmed} \\ 0 & \text{If } \mathbf{x}_i \text{ is trimmed} \end{cases} .$

- **Cases** $\rightarrow \mathbf{x}_i = (x_{i1}, \dots, x_{ip})' \in \mathbb{R}^p$ and **Cells** $\rightarrow x_{ij} \in \mathbb{R}$
 - ◇ i denotes a country (or a trader; company;...) for $i = 1, \dots, n$
 - ◇ x_{ij} is the “quantity-value ratio” for country i in the j -th month (or the j -th year; the j -th product;...) for $j = 1, \dots, p$
 - **Casewise trimming:** Trim \mathbf{x}_i cases with (at least one) outlying x_{ij}
- $n = 100 \times p = 4$ data matrix with 2% outlying cells:



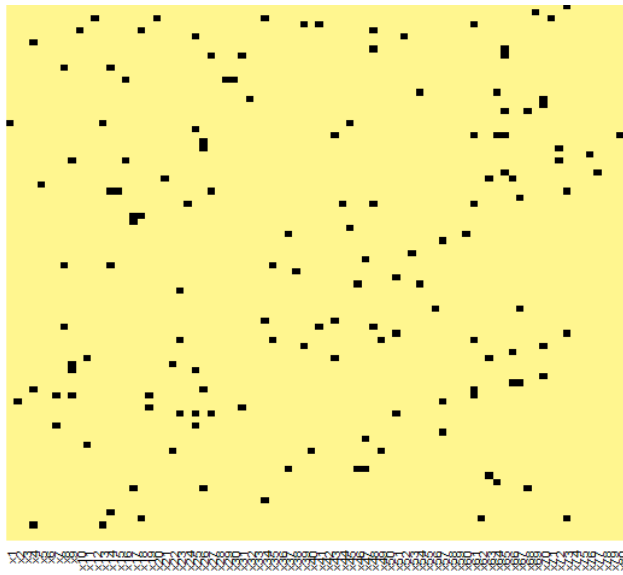
Outlying x_{ij} cells



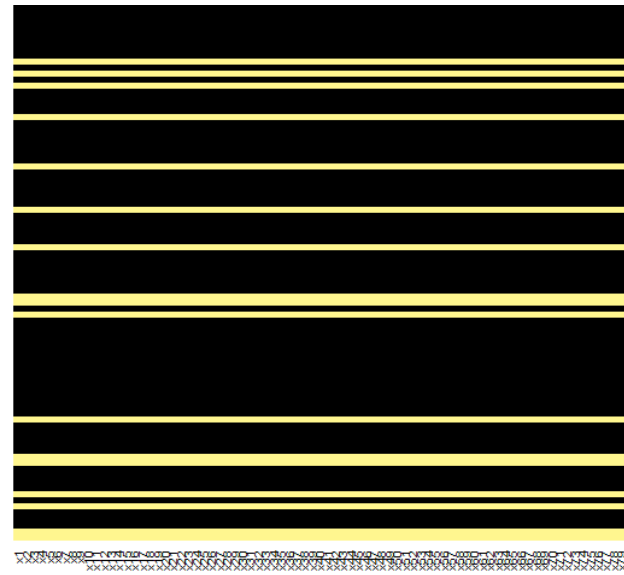
Trimmed \mathbf{x}_i cases (black lines)

- But when the dimension p increases... we do not expect many \mathbf{x}_i completely free of outlying x_{ij} cells:

$n = 100 \times p = 80$ data matrix with 2% outlying cells:



Outlying x_{ij} cells



Trimmed \mathbf{x}_i cases (black lines)

- **Cellwise trimming:**

◇ Only trimming outlying cells... (\Rightarrow "Particular" frauds identified...??)

- PCA approximation $\hat{\mathbf{x}}_i = \mathbf{m} + \mathbf{B}_q \mathbf{a}_i = (\hat{x}_{i1}, \dots, \hat{x}_{ip})^T$ re-written as

$$\hat{x}_{ij} = m_j + \mathbf{a}_i^T \mathbf{b}_j.$$

- **Cellwise LTS** (Cevallos-Valdiviezo 2016): Minimize

$$\sum_{i=1}^n w_{ij} (x_{ij} - m_j - \mathbf{a}_i^T \mathbf{b}_j)^2$$

- ◇ $w_{ij} = 0$ if cell x_{ij} is **trimmed** and $w_{ij} = 1$ if not with

$$\sum_{i=1}^n w_{ij} = [n(1 - \alpha)], \text{ for } j = 1, \dots, p.$$

- Different patterns/structures in data \Rightarrow **G subspace approximations:**

$$\widehat{\mathbf{x}}_i^g(\mathbf{B}_{q_g}^g, \mathbf{A}_{q_g}^g, \mathbf{m}^g) = \mathbf{m}^g + \mathbf{B}_{q_g}^g \mathbf{a}_i^g \quad \text{or} \quad \widehat{x}_{ij}^g = m_j^g + (\mathbf{a}_i^g)^T \mathbf{b}_j^g,$$

for $g = 1, \dots, G$

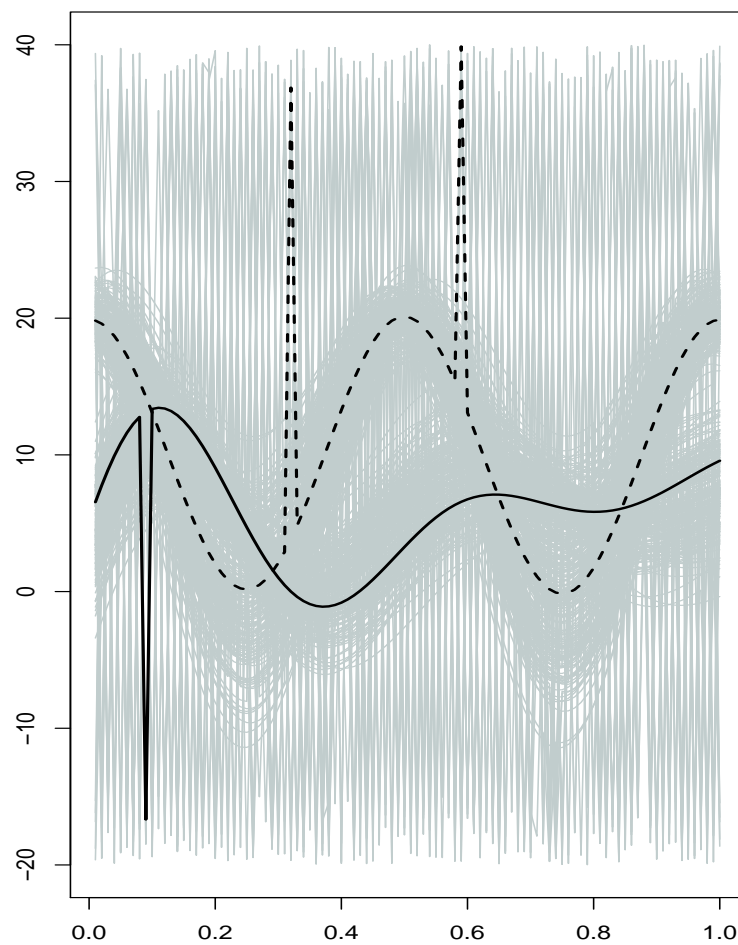
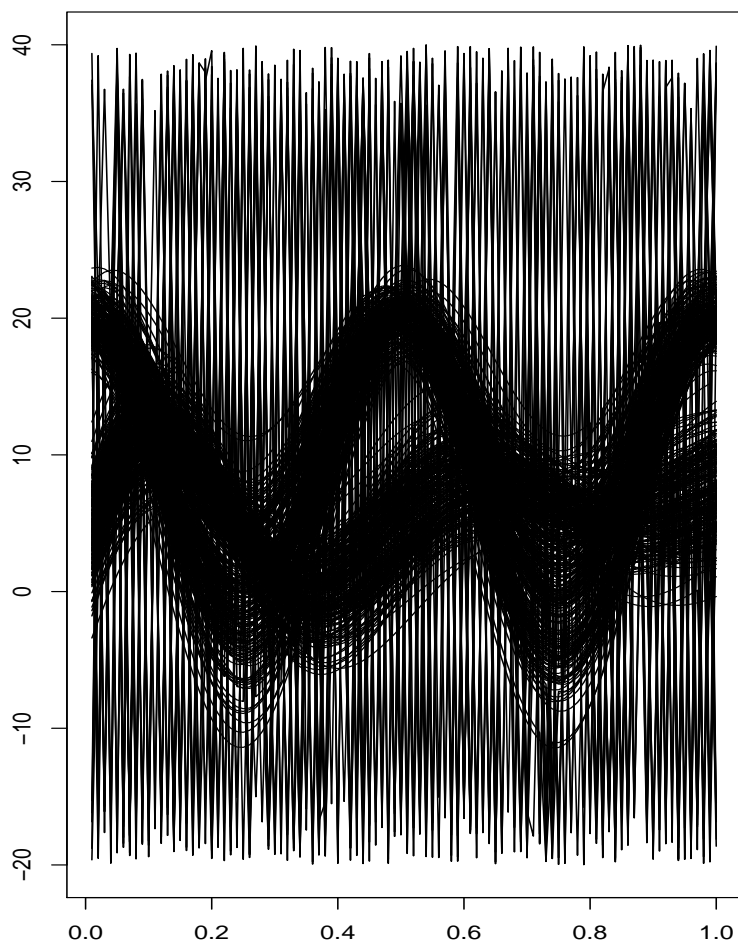
- Minimize

$$\min_{w_{ij}^g, \mathbf{B}_{q_g}^g, \mathbf{A}_{q_g}^g, \mathbf{m}^g} \sum_{i=1}^n \sum_{j=1}^p \sum_{g=1}^G w_{ij}^g (x_{ij} - \widehat{x}_{ij}^g)^2.$$

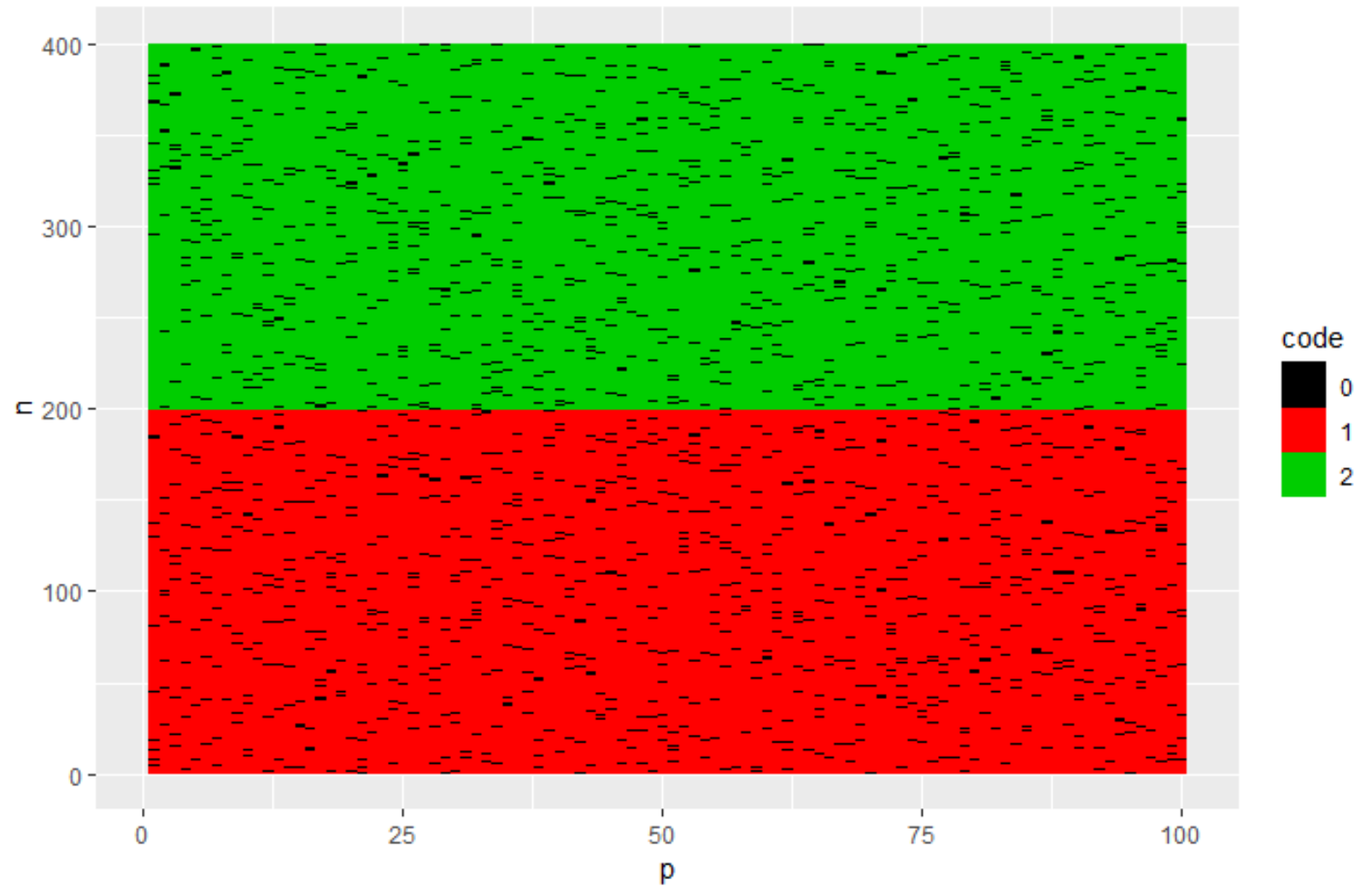
- ◇ $w_{ij}^g = 1$ if cell x_{ij} is assigned to cluster g and non-trimmed and 0 otherwise
- ◇ Appropriate constraints on the w_{ij}^g

q_1, \dots, q_G are intrinsic dimensions...

- **Example 1:** $n = 400$ in dimension $p = 100$ with 2 groups and 2% “scattered” outliers:

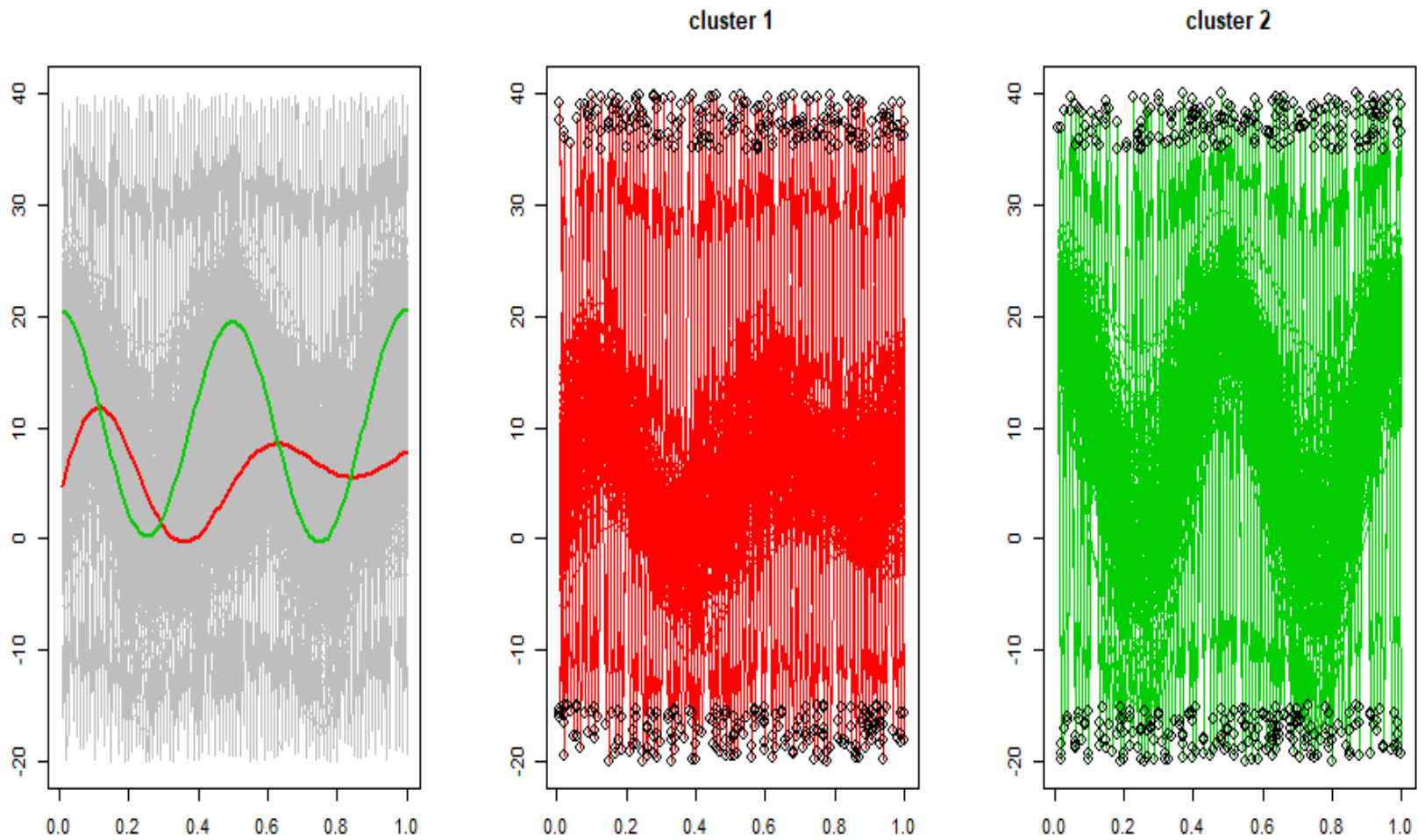


- $k = 2$, $q = 2$ and $\alpha = 0.05$:

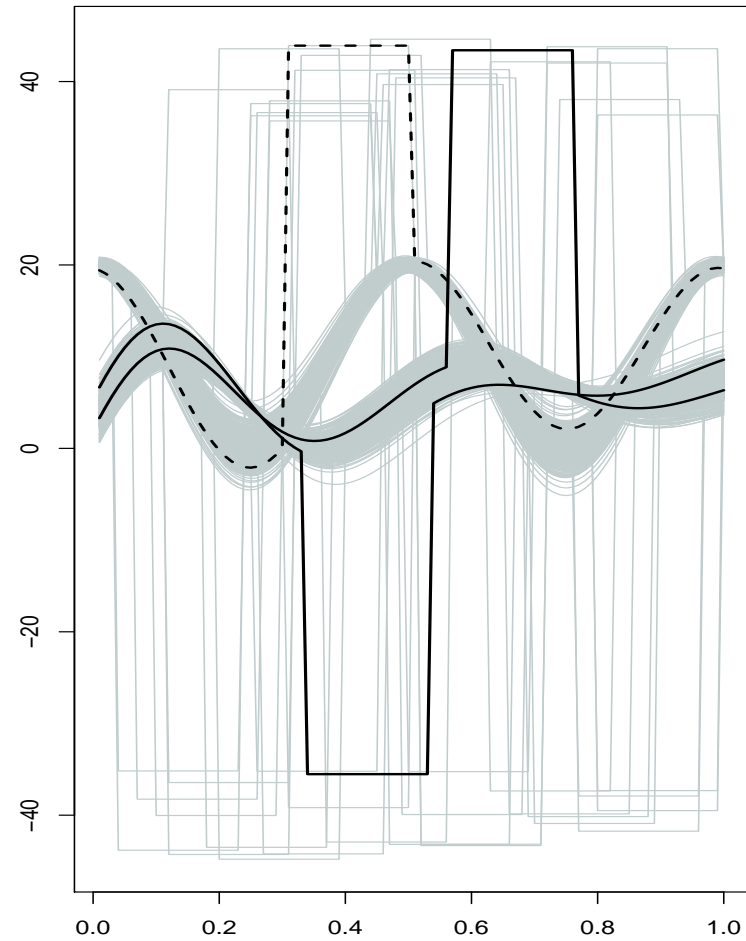
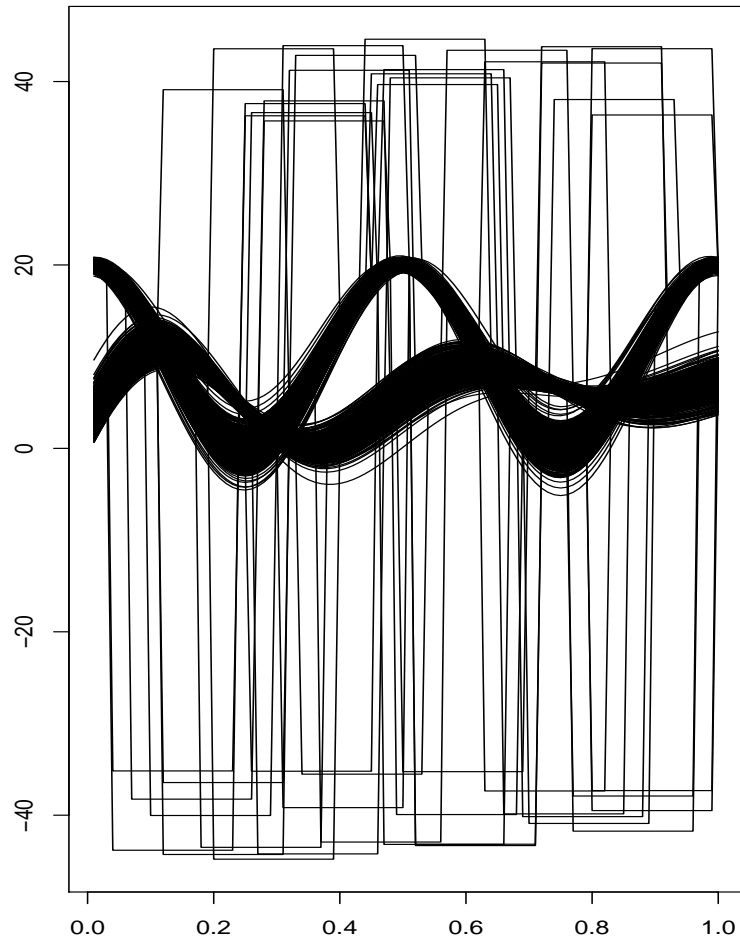


“-” are the trimmed cells

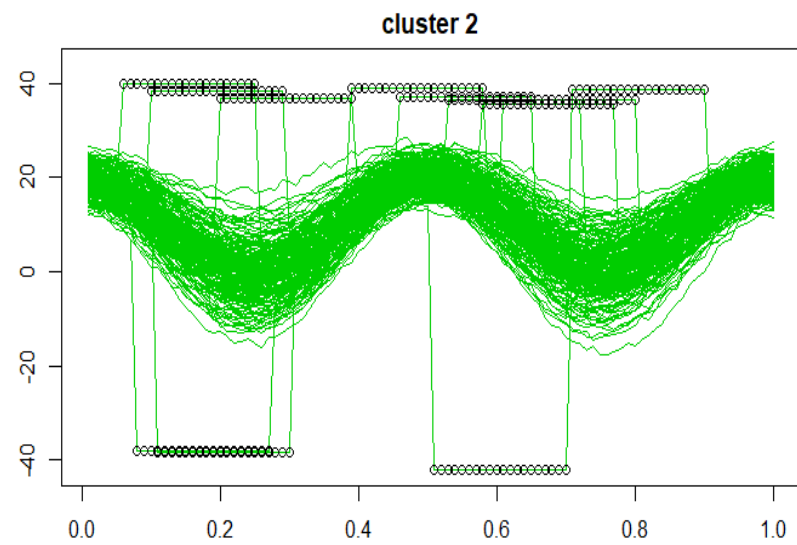
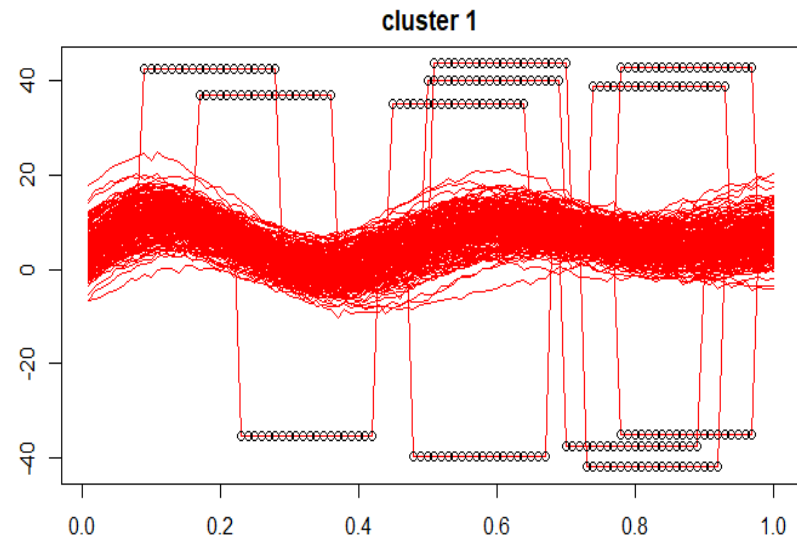
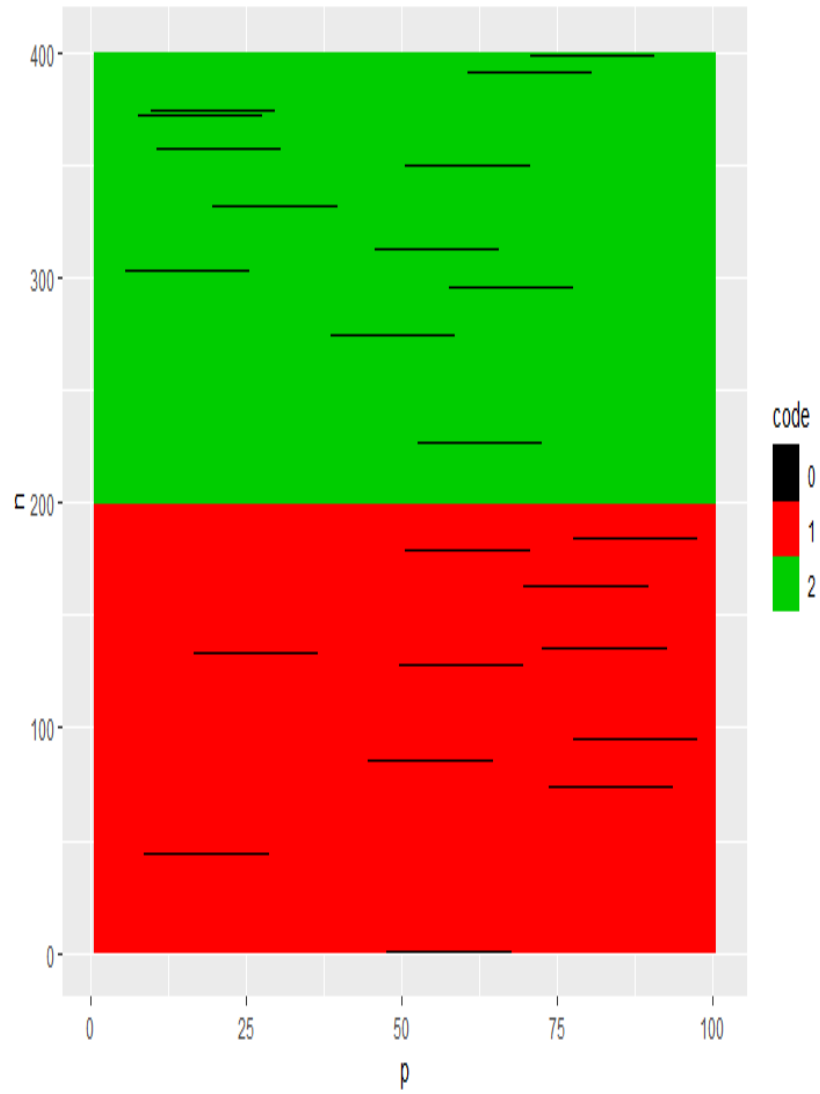
- Cluster means and trimmed cells (\circ):



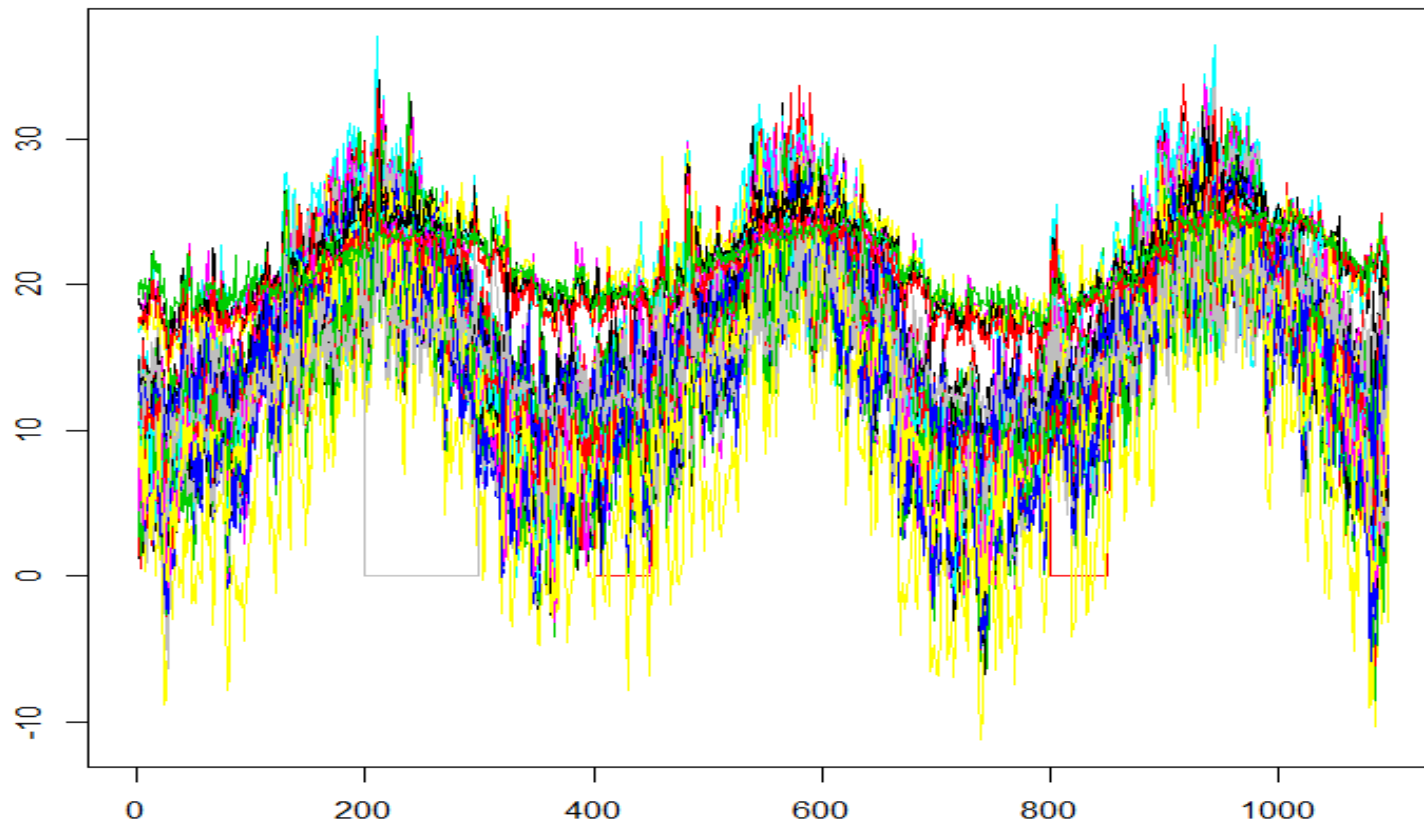
- **Example 2:** $n = 400$ in dimension $p = 100$ with 2 groups and few curves with 20% consecutive cells corrupted:



- Results:

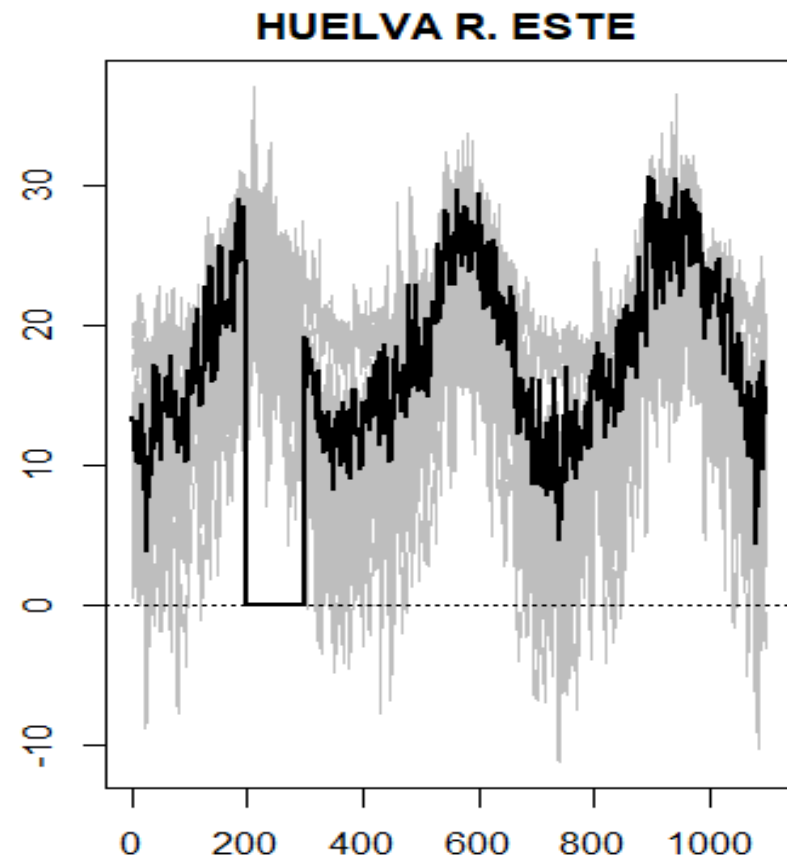
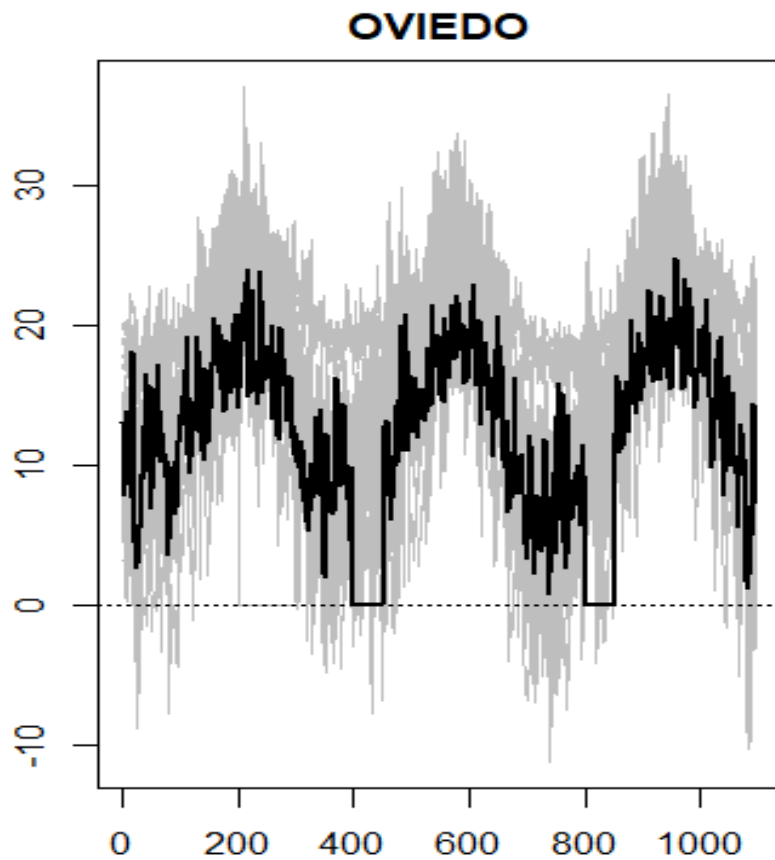


- **Real data example:** Average daily temperatures in 83 Spanish meteorologic stations between 2007-2009 ($n = 83$ and $p = 1096$).



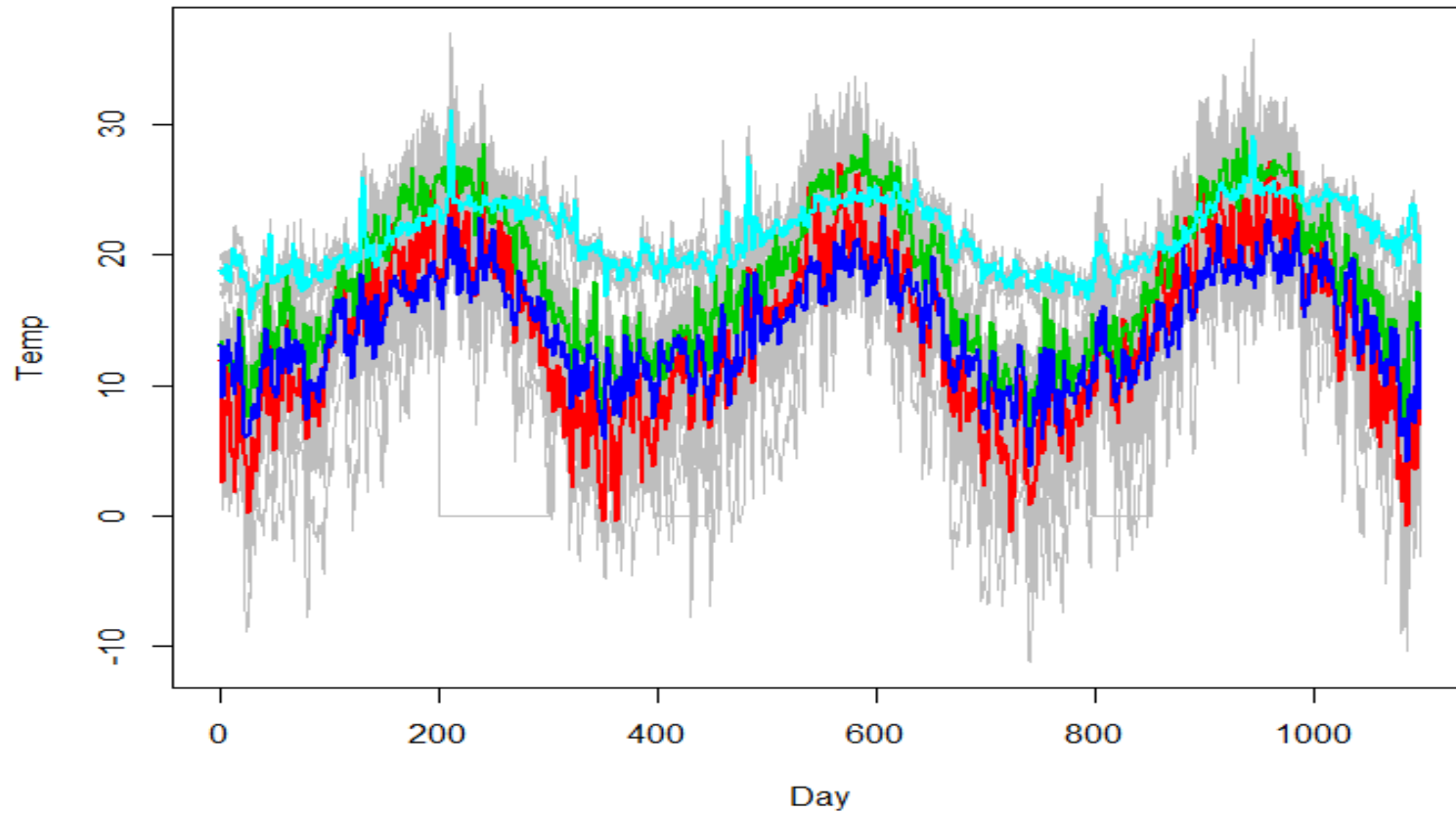
- **Artificial outliers:**

- ◇ Two periods of 50 consecutive days in *Oviedo* replaced by 0°C .
- ◇ 150 consecutive days in *Huelva* temperature replaced by 0°C .

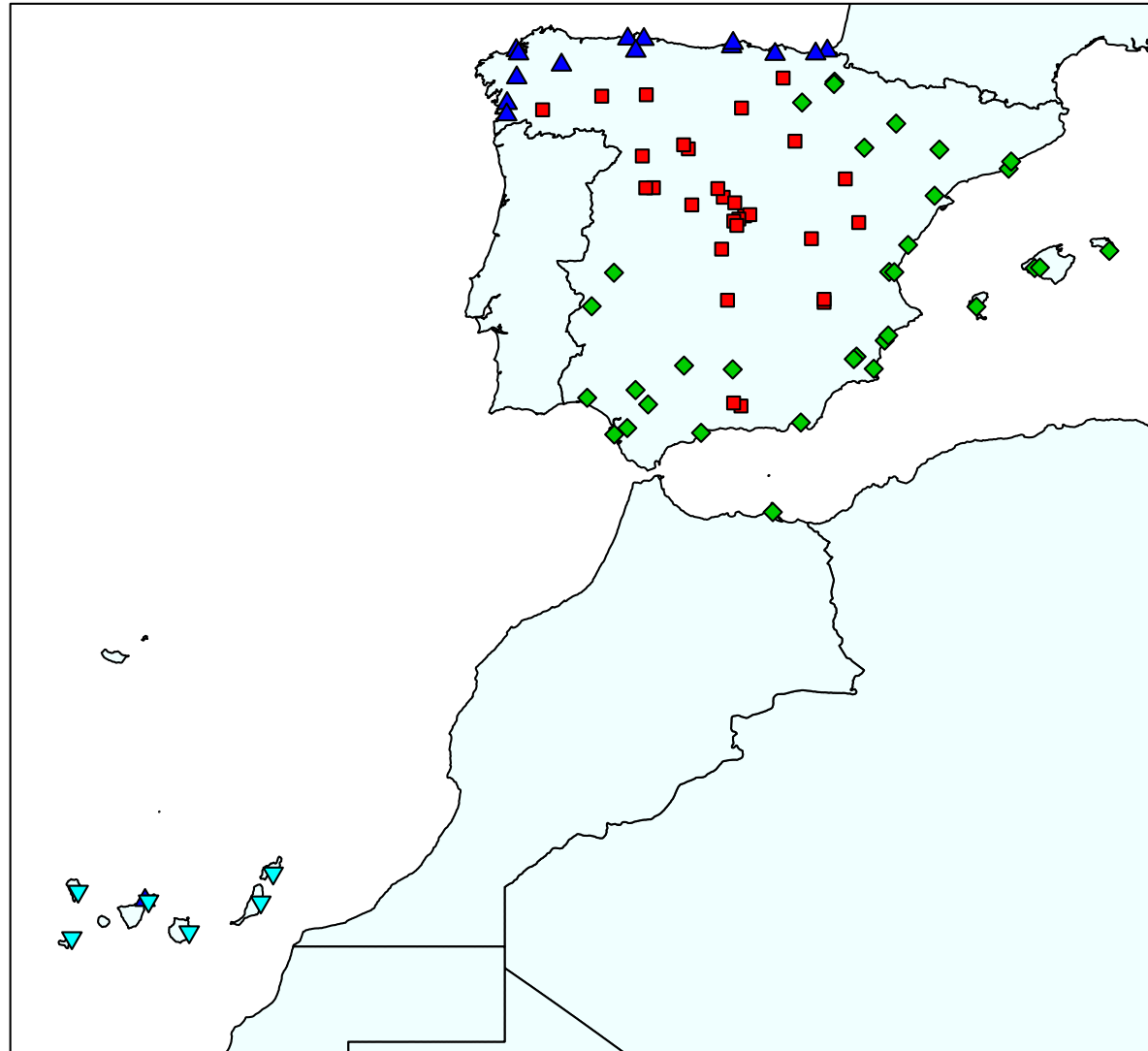


- **Cluster means:**

- ◇ “Meseta” (Central plateau-Castile): — (red)
- ◇ Mediterranean: — (green)
- ◇ Cantabrian Coast: — (blue)
- ◇ Canary Islands: — (cyan)

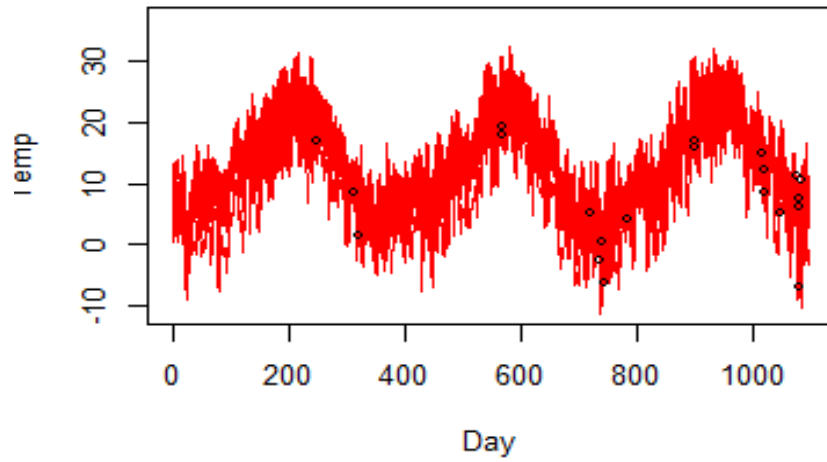


- **Clustered stations:**

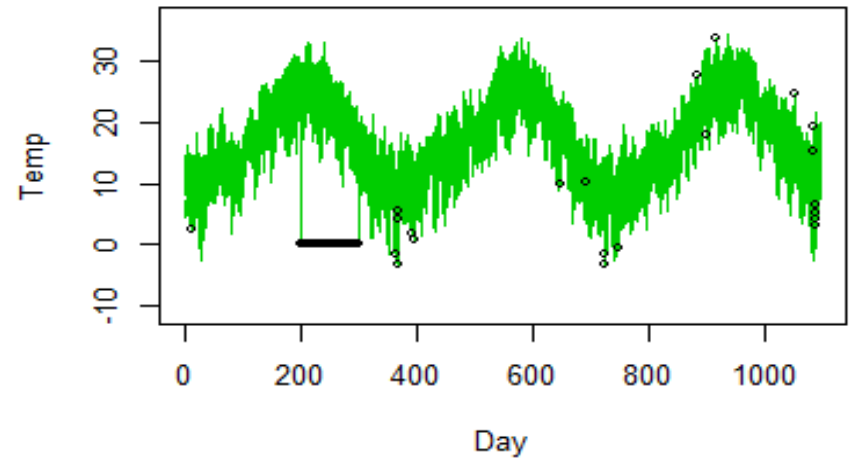


- Clusters found and trimmed cells:

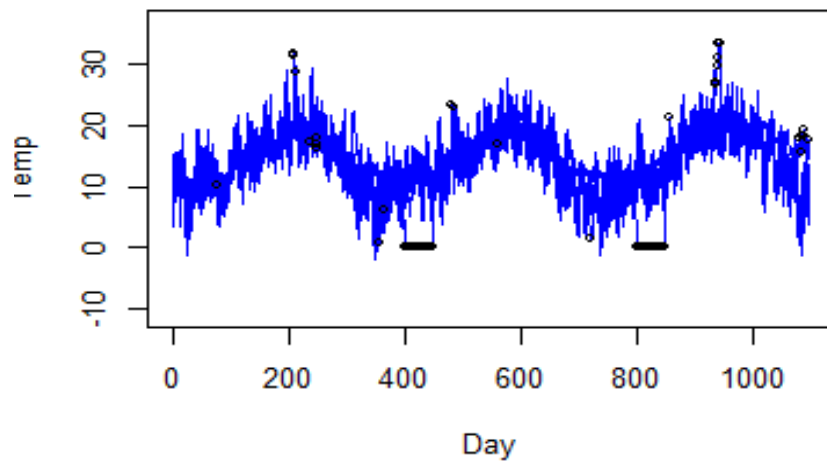
cluster 1



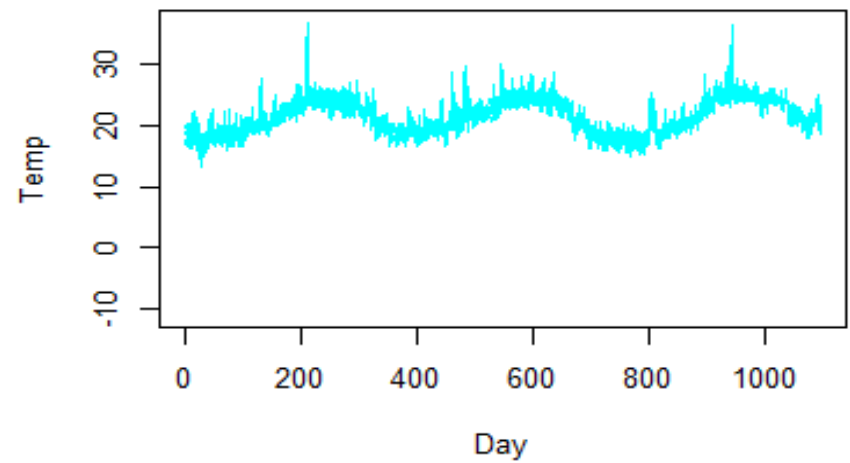
cluster 2



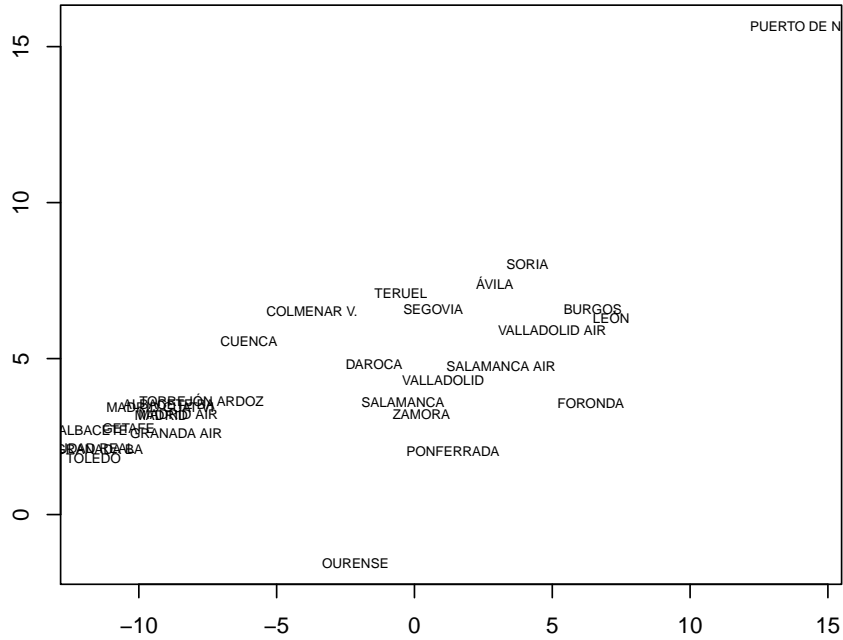
cluster 3



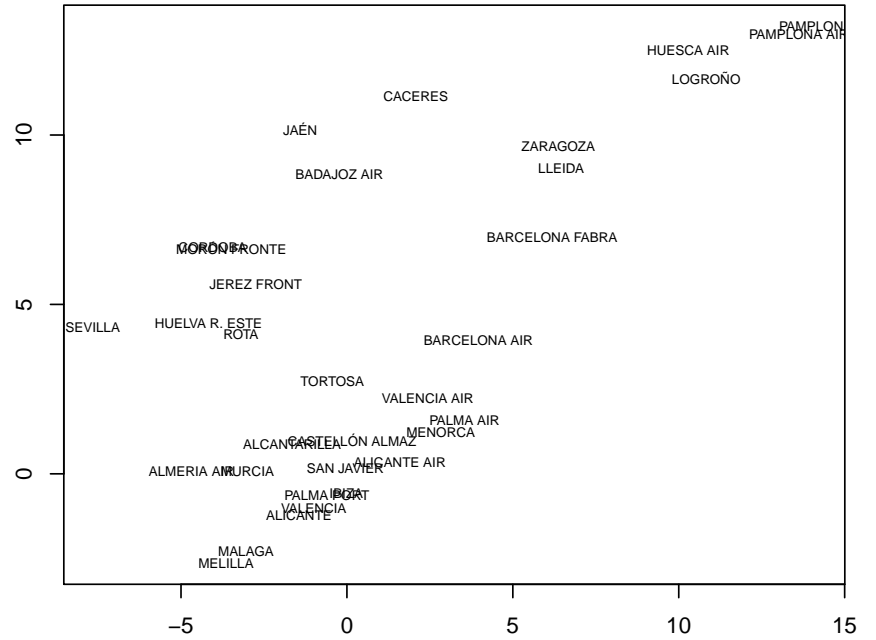
cluster 4



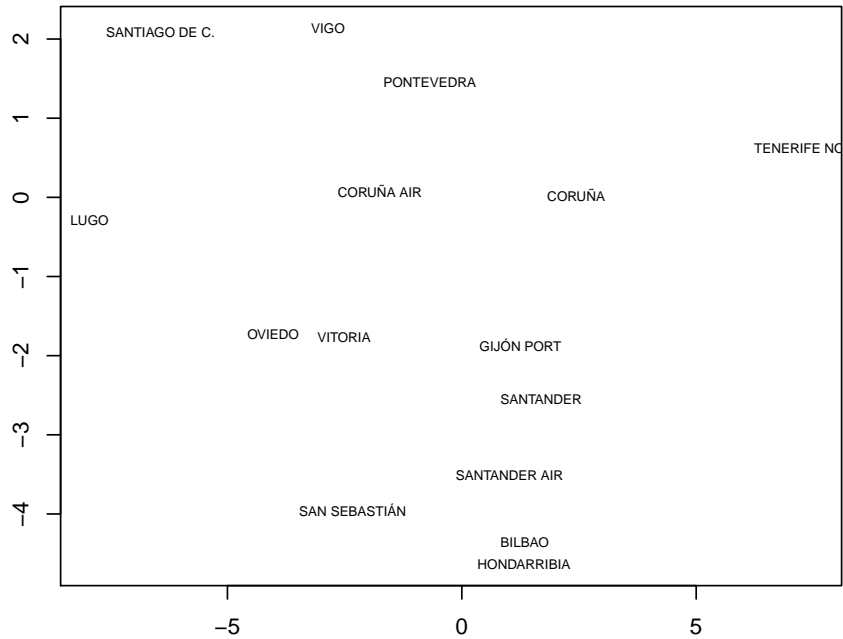
First two scores of cluster 1



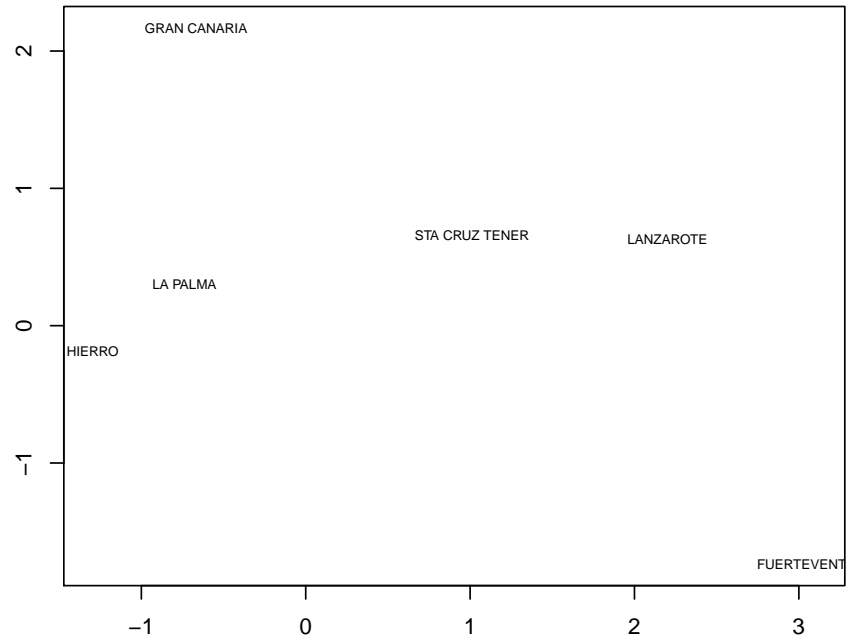
First two scores of cluster 2

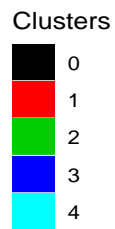
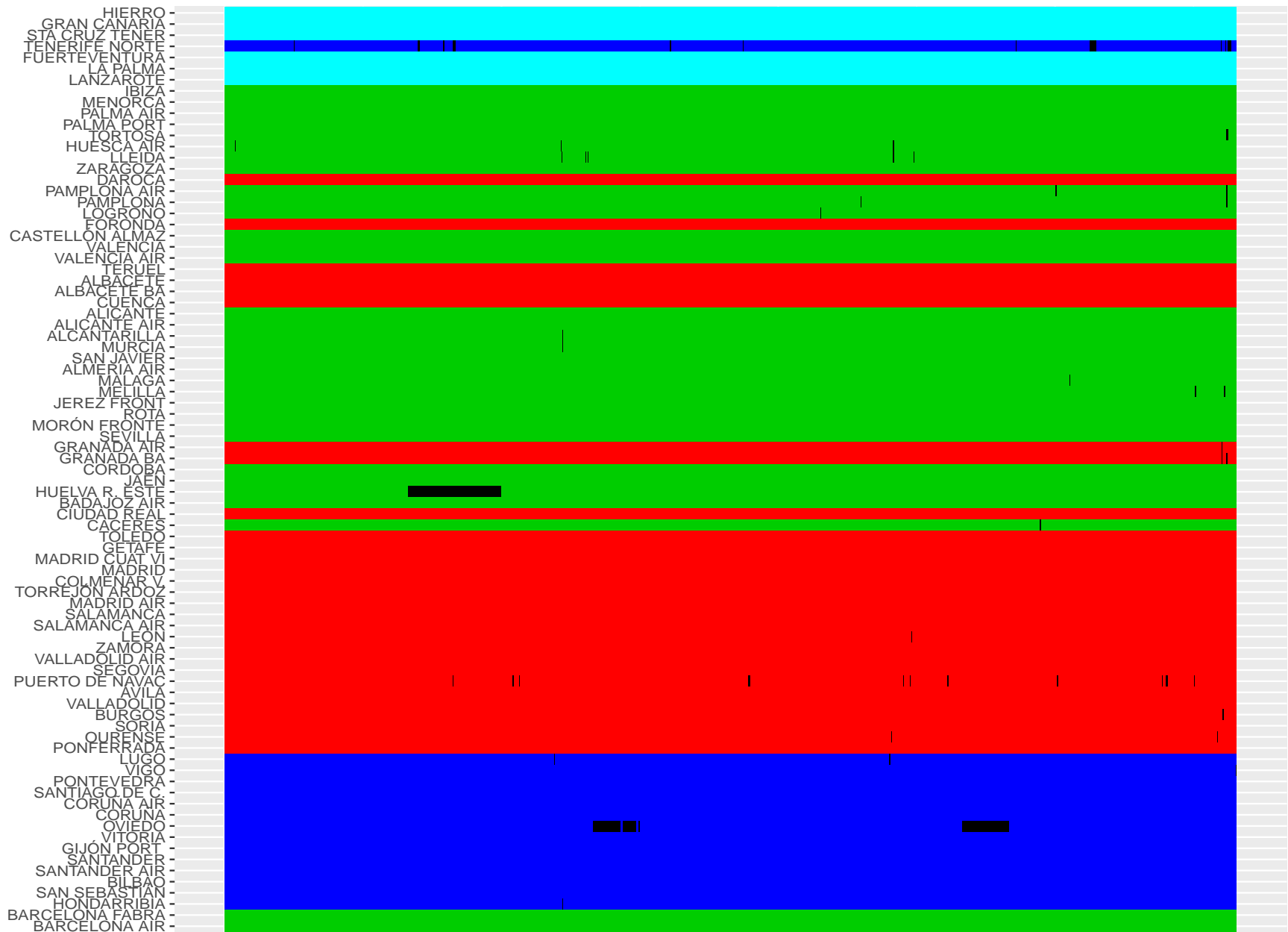


First two scores of cluster 3



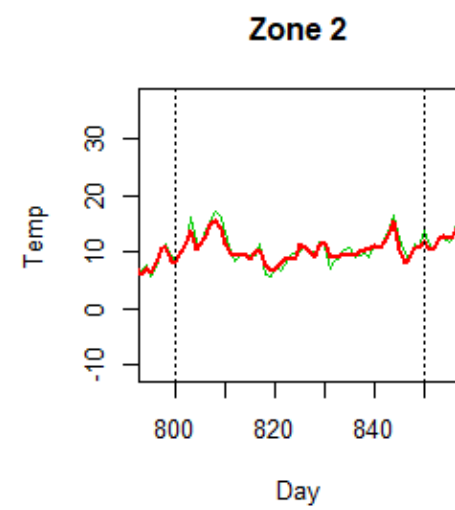
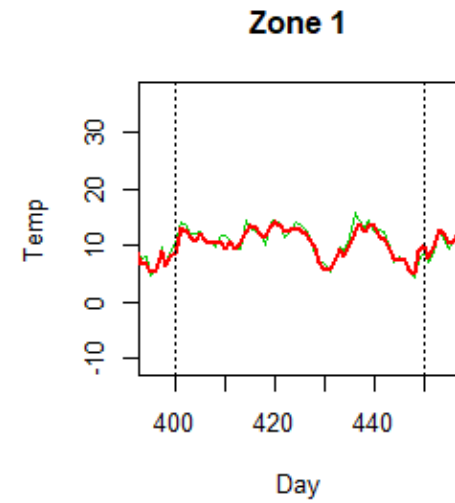
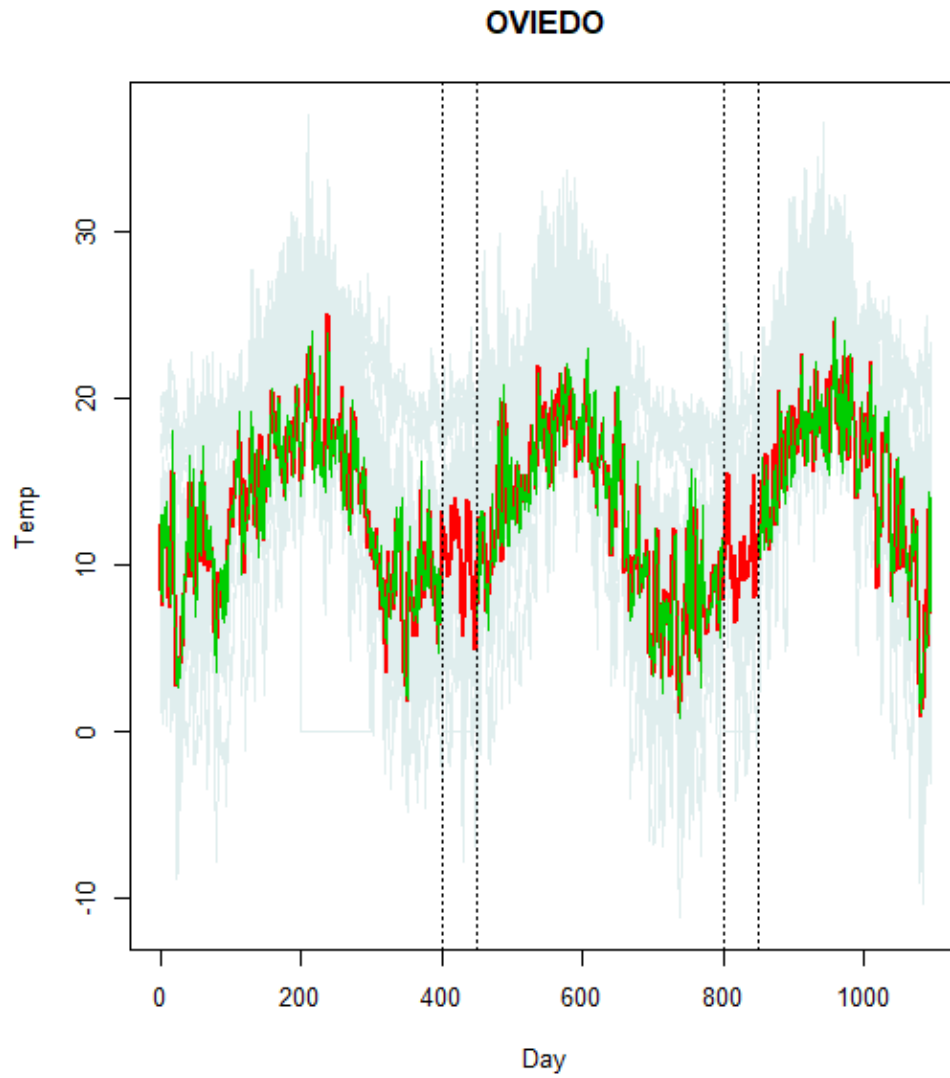
First two scores of cluster 4





Dias

- Reconstructed curves “—” and true real data “—” in Oviedo:



- **Conclusions:**

- ◇ Different patterns/structures in data \Rightarrow Cluster Analysis
- ◇ Robust clustering aimed at (jointly) detecting main clusters (bulk of data) and outliers \Rightarrow Potential “frauds”...
- ◇ Higher dimensional problems: Assume clusters “living” in low-dimensional subspaces
- ◇ “Casewise” and “cellwise” trimming

Some References:

- CUESTA-ALBERTOS, J.A., GORDALIZA, A. AND MATRÁN, C. (1997), “Trimmed k -means: An attempt to robustify quantizers,” *Ann. Statist.*, **25**, 553-576.
- GARCÍA-ESCUADERO, L.A. AND GORDALIZA, A. (1999), “Robustness properties of k -means and trimmed k -means,” *J. Amer. Statist. Assoc.*, **94**, 956-969.
- GARCÍA-ESCUADERO, L.A., GORDALIZA, A., MATRÁN, C. AND MAYO-ISCAR, A. (2008), “A General Trimming Approach To Robust Cluster Analysis,” *Ann. Statist.*, **36**, 1324-1345.
- GARCÍA-ESCUADERO, L.A., GORDALIZA, A., MATRÁN, C. AND MAYO-ISCAR, A. (2010), “A review of robust clustering methods,” *Advances in Data Analysis and Classification*, **4**, 89-109.
- FRITZ, H., GARCÍA-ESCUADERO, L.A. AND MAYO-ISCAR; A (2012), “tclust: An R package for a trimming approach to Cluster Analysis,” *Journal of Statistical Software*, **47**, Issue 12.

Thanks for your attention!!!