

# Robust clustering in (moderately) high dimensional cases

International Conference on Robust Statistics 2025, Stresa,  
Italy

**Luis Angel García-Escudero**

Universidad de Valladolid

*(joint work with Agustín Mayo-Iscar and Lucía Trapote)*

Robust  
clustering in  
higher  
dimensions

L.A. García-  
Escudero

Robust  
clustering and  
trimming

TCLUST high  
dimensions

RLG method

Trimmed  
HDDC

Other issues

Conclusions

- 1 Robust clustering based on trimming
- 2 TCLUST in higher dimensional cases
- 3 Robust Linear Grouping method
- 4 Trimmed HDDC
- 5 Other issues
- 6 Conclusions

# Robust clustering based on trimming



# Why robust clustering?

- Outliers are known to be **problematic** in Cluster Analysis:
  - Relevant and well-defined clusters incorrectly merged
  - Spurious clusters detected

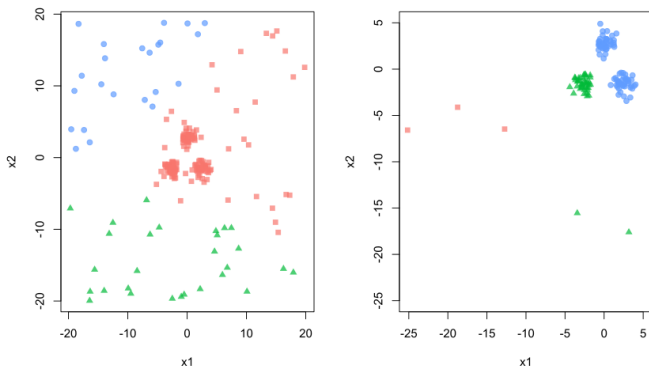


Figure: Impact of contamination on 3-means

- Outliers can be considered as clusters in themselves, suggesting an increase in the number of clusters  $G$ 
  - ◇ Not always the best strategy (and sometimes infeasible...)

# Why robust clustering?

- Outliers can be considered as clusters in themselves, suggesting an increase in the number of clusters  $G$ 
  - ◇ Not always the best strategy (and sometimes infeasible...)
- **Cluster Analysis/Anomaly Detection** are **related** topics:
  - ◇ The first finds crowds of data points, while the second aims to detect observations far from these crowds
  - ◇ An unified treatment: **Robust Clustering**

# Why robust clustering?

- Outliers can be considered as clusters in themselves, suggesting an increase in the number of clusters  $G$ 
  - ◇ Not always the best strategy (and sometimes infeasible...)
- **Cluster Analysis/Anomaly Detection** are **related** topics:
  - ◇ The first finds crowds of data points, while the second aims to detect observations far from these crowds
  - ◇ An unified treatment: **Robust Clustering**
- Clustered outliers are harmful to (even robust) statistical methods but can be easily detected through clustering

- Different approaches for robust clustering:



- Different approaches for robust clustering:
  - ① Noise and outliers accommodated by **heavy-tailed components**: mixtures of  $t$ -distributions [Peel and McLachlan, 2000] or mixtures of contaminated Gaussian distributions [Punzo and McNicholas 2016]

- Different approaches for robust clustering:
  - ① Noise and outliers accommodated by **heavy-tailed components**: mixtures of  $t$ -distributions [Peel and McLachlan, 2000] or mixtures of contaminated Gaussian distributions [Punzo and McNicholas 2016]
  - ② Modelled by a **uniformly distributed noise** component [Banfield and Raftery 1993; Coretto and Hennig 2016]

- Different approaches for robust clustering:
  - ① Noise and outliers accommodated by **heavy-tailed components**: mixtures of  $t$ -distributions [Peel and McLachlan, 2000] or mixtures of contaminated Gaussian distributions [Punzo and McNicholas 2016]
  - ② Modelled by a **uniformly distributed noise** component [Banfield and Raftery 1993; Coretto and Hennig 2016]
  - ③ **Trimming** approach [Cuesta-Albertos et al 1997, Neykov et al 2007, García-Escudero et al 2008]

- Different approaches for robust clustering:
  - ① Noise and outliers accommodated by **heavy-tailed components**: mixtures of  $t$ -distributions [Peel and McLachlan, 2000] or mixtures of contaminated Gaussian distributions [Punzo and McNicholas 2016]
  - ② Modelled by a **uniformly distributed noise** component [Banfield and Raftery 1993; Coretto and Hennig 2016]
  - ③ **Trimming** approach [Cuesta-Albertos et al 1997, Neykov et al 2007, García-Escudero et al 2008]
- We **focus** on that robust clustering approach **based on trimming**

- Standard statistical tools are applied in **trimming** after (hopefully) discarding outliers within the fraction  $\alpha$  of trimmed observations  $\rightsquigarrow$  Easy interpretation

- Standard statistical tools are applied in **trimming** after (hopefully) discarding outliers within the fraction  $\alpha$  of trimmed observations  $\rightsquigarrow$  Easy interpretation
- Trimming **self-determined** by data [Rousseeuw 1984, 1985]:
  - ◇ LTS and LMS in regression
  - ◇ MVE and MCD in location-and-scatter estimation

- Standard statistical tools are applied in **trimming** after (hopefully) discarding outliers within the fraction  $\alpha$  of trimmed observations  $\rightsquigarrow$  Easy interpretation
- Trimming **self-determined** by data [Rousseeuw 1984, 1985]:
  - ◇ LTS and LMS in regression
  - ◇ MVE and MCD in location-and-scatter estimation
- The use of **C-steps (concentration steps)** [Rousseeuw and van Driessen 1999] forms the basis for its practical application

## ADVANCED REVIEW



WILEY

## Robust clustering based on trimming

Luis A. García-Escudero | Agustín Mayo-Iscar

Department of Statistics and Operation  
Research and IMUVA, University of  
Valladolid, Valladolid, Spain

### Correspondence

Luis A. García-Escudero, Department of  
Statistics and Operation Research and  
IMUVA, University of Valladolid,

### Abstract

Clustering is one of the most widely used unsupervised learning techniques. However, it is well-known that outliers can have a significantly adverse impact on commonly applied clustering methods. On the other hand, clustered outliers can be particularly detrimental to (even robust) statistical procedures. Therefore, it makes sense to combine concepts from Robust Statistics and Clus-

Figure: WIREs Comp Stat 2024



# Trimmed $K$ -means [Cuesta-Albertos, Gordaliza and Matrán 1997]

- Given  $\{x_1, \dots, x_n\} \subset \mathbb{R}^p$ , the **trimmed  $k$ -means** search for  $G$  centers  $\mu_1, \dots, \mu_G$  and a partition

$$\{R_0, R_1, \dots, R_G\} \text{ of } \{1, 2, \dots, n\}$$

with

$$\#R_0 = \lfloor n\alpha \rfloor$$

minimizing

$$\sum_{g=1}^G \sum_{i \in R_g} \|x_i - \mu_g\|^2$$

# Trimmed $K$ -means [Cuesta-Albertos, Gordaliza and Matrán 1997]

- Given  $\{x_1, \dots, x_n\} \subset \mathbb{R}^p$ , the **trimmed  $k$ -means** search for  $G$  centers  $\mu_1, \dots, \mu_G$  and a partition

$$\{R_0, R_1, \dots, R_G\} \text{ of } \{1, 2, \dots, n\}$$

with

$$\#R_0 = \lfloor n\alpha \rfloor$$

minimizing

$$\sum_{g=1}^G \sum_{i \in R_g} \|x_i - \mu_g\|^2$$

- $R_1, \dots, R_G$  gives the partition into  $G$  clusters, but a fraction  $\alpha$  of observations (those with indexes in  $R_0$ ) are trimmed

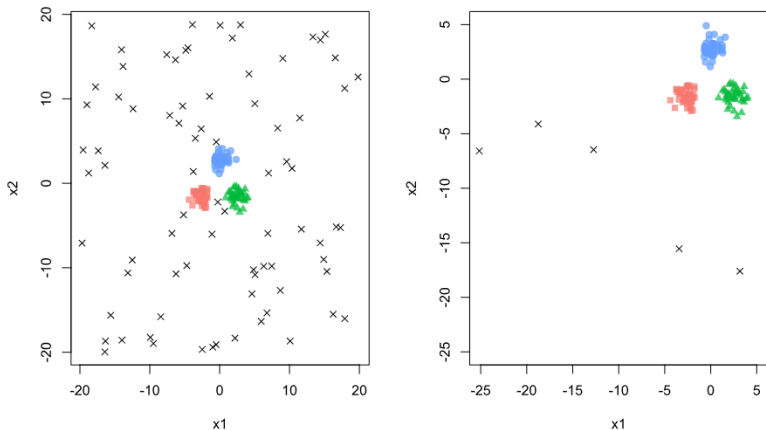


Figure: Trimmed 3-means with  $\alpha = 0.05$ . Trimmed points as “x”

# Trimmed $K$ -means alg. [G-E, Gordaliza and Matrán 2003]

## 1 *Random initializations:*

$G$  random observations  $\rightsquigarrow$  Initial centers  $\mu_1^0, \dots, \mu_G^0$

Trimmed  $K$ -means alg. [G-E, Gordaliza and Matrán 2003]① *Random initializations:* $G$  random observations  $\rightsquigarrow$  Initial centers  $\mu_1^0, \dots, \mu_G^0$ ② *C-steps:*

## 2.1 Take

$$D_{ig} = \|x_i - \mu_g^{t-1}\|^2 \text{ and } D_i = \min_{g=1, \dots, G} D_{ig}$$

and  $D_{(1)} \leq D_{(2)} \leq \dots \leq D_{(n)}$  to get

$$R_g = \{i : D_{ig} = D_i \text{ and } D_i \leq D_{([n(1-\alpha)])}\}$$

for  $g = 1, \dots, G$ 2.2 Update centers  $\mu_g^t = \mathbf{avg}\{x_i : i \in R_g\}$

Trimmed  $K$ -means alg. [G-E, Gordaliza and Matrán 2003]① *Random initializations:* $G$  random observations  $\rightsquigarrow$  Initial centers  $\mu_1^0, \dots, \mu_G^0$ ② *C-steps:*

## 2.1 Take

$$D_{ig} = \|x_i - \mu_g^{t-1}\|^2 \text{ and } D_i = \min_{g=1, \dots, G} D_{ig}$$

and  $D_{(1)} \leq D_{(2)} \leq \dots \leq D_{(n)}$  to get

$$R_g = \{i : D_{ig} = D_i \text{ and } D_i \leq D_{([n(1-\alpha)])}\}$$

for  $g = 1, \dots, G$ 2.2 Update centers  $\mu_g^t = \mathbf{avg}\{x_i : i \in R_g\}$ ③ *Output:* that with the minimum value of the target function.

# Trimmed $K$ -means alg. [G-E, Gordaliza and Matrán 2003]

## 1 *Random initializations:*

$G$  random observations  $\rightsquigarrow$  Initial centers  $\mu_1^0, \dots, \mu_G^0$

## 2 *C-steps:*

### 2.1 Take

$$D_{ig} = \|x_i - \mu_g^{t-1}\|^2 \text{ and } D_i = \min_{g=1, \dots, G} D_{ig}$$

and  $D_{(1)} \leq D_{(2)} \leq \dots \leq D_{(n)}$  to get

$$R_g = \{i : D_{ig} = D_i \text{ and } D_i \leq D_{([n(1-\alpha)])}\}$$

for  $g = 1, \dots, G$

### 2.2 Update centers $\mu_g^t = \mathbf{avg}\{x_i : i \in R_g\}$

## 3 *Output:* that with the minimum value of the target function.

- It reduces to Lloyd's  $k$ -means algorithm when  $\alpha = 0$

# TCLUST as an extension of trimmed $k$ -means

- Trimmed  $k$ -means favours spherical/equally scattered clusters

Robust  
clustering and  
trimming

TCLUST high  
dimensions

RLG method

Trimmed  
HDDC

Other issues

Conclusions



# TCLUST as an extension of trimmed $k$ -means

- Trimmed  $k$ -means favours spherical/equally scattered clusters
- More **flexible** clustering associated to  $G$  **normal** components with location  $\{\mu_g\}_{g=1}^G$  and **scatter** matrices  $\{\Sigma_g\}_{g=1}^G$

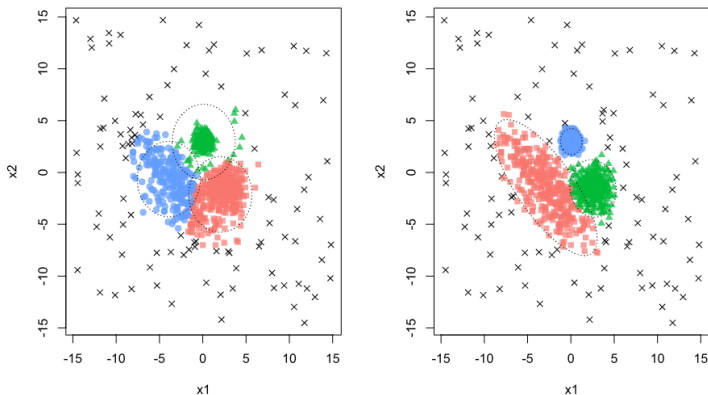


Figure: Trimmed 3-means (left). TCLUST  $G = 3, c = 12$  and  $\alpha = 0.05$  (right)

- **TCLUST** searches for centers  $\mu_1, \dots, \mu_G$ , **scatter matrices**  $\Sigma_1, \dots, \Sigma_G$ , **weights**  $\pi_1, \dots, \pi_g$  (with  $\sum_{g=1}^G \pi_g = 1$ ), and a partition  $\{R_0, R_1, \dots, R_G\}$  with  $\#R_0 = [n\alpha]$  maximizing

$$\sum_{g=1}^G \sum_{i \in R_g} \log [\pi_g \phi(x_i; \mu_g, \Sigma_g)],$$

where  $\phi(\cdot; \mu, \Sigma)$  is the *p.d.f.* of a *p*-variate normal

- Another important ingredient of TCLUST is the **eigenvalues-ratio constraint**:

$$\frac{\max_{g=1,\dots,G,j=1,\dots,p} \lambda_j(\Sigma_g)}{\min_{g=1,\dots,G,j=1,\dots,p} \lambda_j(\Sigma_g)} \leq c,$$

where  $\{\lambda_j(\Sigma)\}_{j=1}^p$  is the set of eigenvalues of  $\Sigma$  and  $c \geq 1$

- Another important ingredient of TCLUST is the **eigenvalues-ratio constraint**:

$$\frac{\max_{g=1,\dots,G,j=1,\dots,p} \lambda_j(\Sigma_g)}{\min_{g=1,\dots,G,j=1,\dots,p} \lambda_j(\Sigma_g)} \leq c,$$

where  $\{\lambda_j(\Sigma)\}_{j=1}^p$  is the set of eigenvalues of  $\Sigma$  and  $c \geq 1$

- Trimmed  $K$ -means when  $c = 1$  (and  $\pi_1 = \dots = \pi_G$ )

- Another important ingredient of TCLUST is the **eigenvalues-ratio constraint**:

$$\frac{\max_{g=1,\dots,G,j=1,\dots,p} \lambda_j(\Sigma_g)}{\min_{g=1,\dots,G,j=1,\dots,p} \lambda_j(\Sigma_g)} \leq c,$$

where  $\{\lambda_j(\Sigma)\}_{j=1}^p$  is the set of eigenvalues of  $\Sigma$  and  $c \geq 1$

- Trimmed  $K$ -means when  $c = 1$  (and  $\pi_1 = \dots = \pi_G$ )
- Any  $c < \infty$  makes the constrained maximization of the trimmed likelihood well-defined (target function unbounded when  $\mu_g = x_i$  and  $|\Sigma_g| \downarrow 0$ )

- Another important ingredient of TCLUST is the **eigenvalues-ratio constraint**:

$$\frac{\max_{g=1,\dots,G,j=1,\dots,p} \lambda_j(\Sigma_g)}{\min_{g=1,\dots,G,j=1,\dots,p} \lambda_j(\Sigma_g)} \leq c,$$

where  $\{\lambda_j(\Sigma)\}_{j=1}^p$  is the set of eigenvalues of  $\Sigma$  and  $c \geq 1$

- Trimmed  $K$ -means when  $c = 1$  (and  $\pi_1 = \dots = \pi_G$ )
- Any  $c < \infty$  makes the constrained maximization of the trimmed likelihood well-defined (target function unbounded when  $\mu_g = x_i$  and  $|\Sigma_g| \downarrow 0$ )
- Also prevents detecting “spurious” clusters

- A different, but also relevant, source of **lack of robustness**
- **Spurious clusters** (non-interesting clusters formed by a few almost collinear observations):

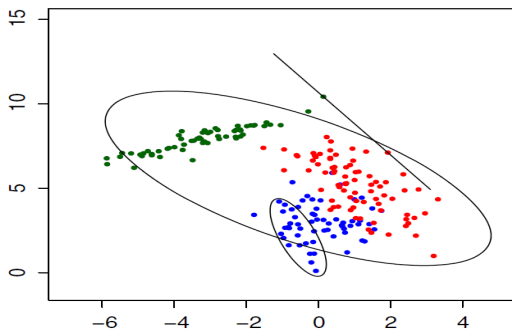


Figure: A spurious cluster detected with  $|\Sigma_g| \approx 0$

- A different, but also relevant, source of **lack of robustness**
- **Spurious clusters** (non-interesting clusters formed by a few almost collinear observations):

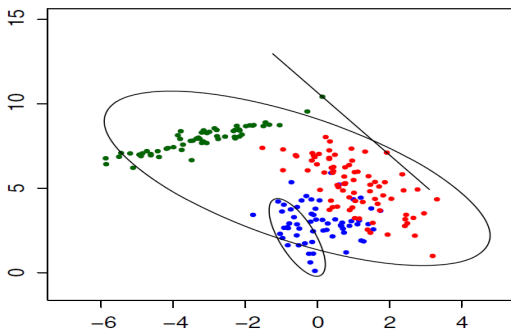


Figure: A spurious cluster detected with  $|\Sigma_g| \approx 0$



- MCD in the case  $G = 1$  (for a large  $c$ )

- MCD in the case  $G = 1$  (for a large  $c$ )
- **Trimmed mixture likelihoods** [Neykov, Filzmoser, Dimova and Neytchev 2007;  $G$ -E, Gordaliza and Mayo-Iscar 2014] maximizing

$$\sum_{i \in R} \log \left[ \sum_{g=1}^G p_g \phi(x_i; m_g, S_g) \right]$$

with  $\#R = [n(1 - \alpha)]$

## 1 *Random initializations:*

$G \times (p + 1)$  random observations  $\rightsquigarrow \mu_1^0, \dots, \mu_G^0$  and  $\Sigma_1^0, \dots, \Sigma_G^0$

① *Random initializations:*

$G \times (p + 1)$  random observations  $\rightsquigarrow \mu_1^0, \dots, \mu_G^0$  and  $\Sigma_1^0, \dots, \Sigma_G^0$

② *C-steps:*

## 2.1 Take

$$D_{ig} = \pi_g^{t-1} \phi(x_i, \mu_g^{t-1}, \Sigma_g^{t-1}) \text{ and } D_i = \max_{g=1, \dots, G} D_{ig}$$

and  $D_{(1)} \leq D_{(2)} \leq \dots \leq D_{(n)}$  to get

$$R_g = \{i : D_{ig} = D_i \text{ and } D_i \geq D_{([n\alpha])}\}$$

① *Random initializations:*

$G \times (p + 1)$  random observations  $\rightsquigarrow \mu_1^0, \dots, \mu_G^0$  and  $\Sigma_1^0, \dots, \Sigma_G^0$

② *C-steps:*

## 2.1 Take

$$D_{ig} = \pi_g^{t-1} \phi(x_i, \mu_g^{t-1}, \Sigma_g^{t-1}) \text{ and } D_i = \max_{g=1, \dots, G} D_{ig}$$

and  $D_{(1)} \leq D_{(2)} \leq \dots \leq D_{(n)}$  to get

$$R_g = \{i : D_{ig} = D_i \text{ and } D_i \geq D_{([n\alpha])}\}$$

2.2 **MLE** $\{x_i : i \in R_g\} \rightsquigarrow \pi_g^t, \mu_g^t$  and  $S_g$

## 1 Random initializations:

$G \times (p + 1)$  random observations  $\rightsquigarrow \mu_1^0, \dots, \mu_G^0$  and  $\Sigma_1^0, \dots, \Sigma_G^0$

## 2 C-steps:

### 2.1 Take

$$D_{ig} = \pi_g^{t-1} \phi(x_i, \mu_g^{t-1}, \Sigma_g^{t-1}) \text{ and } D_i = \max_{g=1, \dots, G} D_{ig}$$

and  $D_{(1)} \leq D_{(2)} \leq \dots \leq D_{(n)}$  to get

$$R_g = \{i : D_{ig} = D_i \text{ and } D_i \geq D_{([n\alpha])}\}$$

**2.2 MLE**  $\{x_i : i \in R_g\} \rightsquigarrow \pi_g^t, \mu_g^t$  and  $S_g$

**2.3** Eigenvalue-ratio constraint imposed on the  $S_g$  to update  $\Sigma_g^t$

### 1 *Random initializations:*

$G \times (p + 1)$  random observations  $\rightsquigarrow \mu_1^0, \dots, \mu_G^0$  and  $\Sigma_1^0, \dots, \Sigma_G^0$

### 2 *C-steps:*

#### 2.1 Take

$$D_{ig} = \pi_g^{t-1} \phi(x_i, \mu_g^{t-1}, \Sigma_g^{t-1}) \text{ and } D_i = \max_{g=1, \dots, G} D_{ig}$$

and  $D_{(1)} \leq D_{(2)} \leq \dots \leq D_{(n)}$  to get

$$R_g = \{i : D_{ig} = D_i \text{ and } D_i \geq D_{([n\alpha])}\}$$

2.2 **MLE**  $\{x_i : i \in R_g\} \rightsquigarrow \pi_g^t, \mu_g^t$  and  $S_g$

2.3 Eigenvalue-ratio constraint imposed on the  $S_g$  to update  $\Sigma_g^t$

### 3 *Output* that with the maximum value

- The **optimal truncation operator**:

$$S_g = \mathbf{cov}\{x_i : i \in R_g\} = U_g \text{diag}(d_{g1}, \dots, d_{gp}) U_g',$$

with  $U_g' U_g = \mathbf{I}_p$ .



- The **optimal truncation operator**:

$$S_g = \mathbf{cov}\{x_i : i \in R_g\} = U_g \text{diag}(d_{g1}, \dots, d_{gp}) U_g',$$

with  $U_g' U_g = \mathbf{I}_p$ . If

$$[d_{gl}]_m = \max\{m, \min\{d_{gl}, c \cdot m\}\},$$

- The **optimal truncation operator**:

$$S_g = \mathbf{cov}\{x_i : i \in R_g\} = U_g \text{diag}(d_{g1}, \dots, d_{gp}) U_g',$$

with  $U_g' U_g = \mathbf{I}_p$ . If

$$[d_{gl}]_m = \max\{m, \min\{d_{gl}, c \cdot m\}\},$$

$n_g = \#R_g$ , and

$$m^* = \arg \min_m \sum_{g=1}^G n_g \sum_{l=1}^p \left( \log[d_{gl}]_m + \frac{d_{gl}}{\log[d_{gl}]_m} \right),$$

then the optimal update is

$$\Sigma_g^t = U_g \text{diag}([d_{g1}]_{m^*}, \dots, [d_{gp}]_{m^*}) U_g'$$

- **Population** version for a theoretical  $P$  and its **empirical** version from an i.i.d sample  $\{X_1, \dots, X_n\}$  with  $X_i \sim P$

- **Population** version for a theoretical  $P$  and its **empirical** version from an i.i.d sample  $\{X_1, \dots, X_n\}$  with  $X_i \sim P$
- **Consistency** of empirical toward population solution

- **Population** version for a theoretical  $P$  and its **empirical** version from an i.i.d sample  $\{X_1, \dots, X_n\}$  with  $X_i \sim P$
- **Consistency** of empirical toward population solution
- **Good robustness** performance:
  - ◇ Influence function [Ruwet et al 2012]
  - ◇ Breakdown point [Ruwet et al 2013]

- **Population** version for a theoretical  $P$  and its **empirical** version from an i.i.d sample  $\{X_1, \dots, X_n\}$  with  $X_i \sim P$
- **Consistency** of empirical toward population solution
- **Good robustness** performance:
  - ◇ Influence function [Ruwet et al 2012]
  - ◇ Breakdown point [Ruwet et al 2013]
- An efficient algorithm and **packages** are available:
  - ◇ **tclust** package in R [Fritz, G-E and Mayo-Isacar 2012]
  - ◇ **FSDA** Matlab toolbox [Cerioli, Riani, Atkinson and Corbellini 2018]

# TCLUST in higher dimensional cases



# Higher dimensions and TCLUST...

- **Higher-dimensional** problems are increasingly **common** in current statistical practice



# Higher dimensions and TCLUST...

- **Higher-dimensional** problems are increasingly **common** in current statistical practice
- TCLUST works well for small dimensions, but face **challenges** as  $p$  increases:

- **Higher-dimensional** problems are increasingly **common** in current statistical practice
- TCLUST works well for small dimensions, but face **challenges** as  $p$  increases:
  - ① Difficulties in **initialization**

- **Higher-dimensional** problems are increasingly **common** in current statistical practice
- TCLUST works well for small dimensions, but face **challenges** as  $p$  increases:
  - 1 Difficulties in **initialization**
  - 2 High **number of parameters** involved

- **Digits data:** A sample of  $n = 1756$  handwritten digits (“3”, “5” and “8”) from the US postal services (subset of a dataset from UCI). Each digit is a  $16 \times 16$  gray level image resulting in  $p = 16^2 = 256$  dimensional vectors:

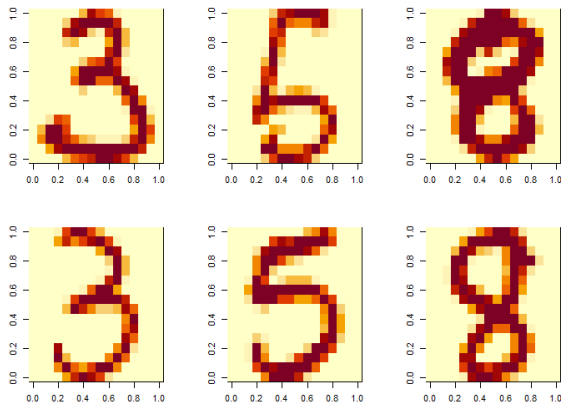


Figure: Six digits images from the Digits dataset

- **TCLUST** applied to the **Digits data** with  $G = 3$ ,  $c = 12$  and  $\alpha = 0.05$  produces (clusters shown in rows, with trimmed images as 0, and actual digits in columns):

	3	5	8
0	21	38	29
1	438	259	81
2	18	61	418
3	181	198	14

- Not satisfactory classification results

- **TCLUST** applied to the **Digits data** with  $G = 3$ ,  $c = 12$  and  $\alpha = 0.05$  produces (clusters shown in rows, with trimmed images as 0, and actual digits in columns):

	3	5	8
0	21	38	29
1	438	259	81
2	18	61	418
3	181	198	14

- Not satisfactory classification results
- Requires considerable computing time

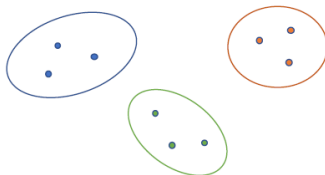
# Difficulties in initialization

- The algorithm **ideally** requires for initialization a random sub-sample of size  $G \times (p+1)$  being outlier-free and properly grouped:

$G \times (p+1)$  observations  $\Rightarrow$

First $p+1$	$\hookrightarrow$ Cluster 1
Second $p+1$	$\hookrightarrow$ Cluster 2
$\vdots$	
$G$ -th $p+1$	$\hookrightarrow$ Cluster $G$

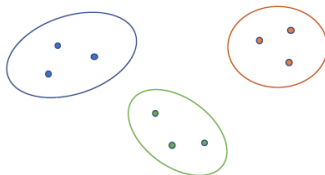
1
2
3
4
5
6
7
8
9



- The algorithm **ideally** requires for initialization a random sub-sample of size  $G \times (p+1)$  being outlier-free and properly grouped:

$G \times (p+1)$ observations	$\Rightarrow$	First $p+1$	$\hookrightarrow$ Cluster 1
		Second $p+1$	$\hookrightarrow$ Cluster 2
		$\vdots$	
		$G$ -th $p+1$	$\hookrightarrow$ Cluster $G$

1
2
3
4
5
6
7
8
9



- That appropriate initializations becomes increasingly **unlikely**



- **Recent** implementation of the **tclust** package (**version 2.0**: jointly with V. Todorov):
  - ① **nstart** random initializations with small **niter1** of C-steps
  - ② **nkeep** solutions with the largest values which are iterated until convergence or further **niter2** C-steps

# Difficulties in initialization: possible remedies

- **Recent** implementation of the **tclust** package (**version 2.0**: jointly with V. Todorov):

- ① **nstart** random initializations with small **niter1** of C-steps
- ② **nkeep** solutions with the largest values which are iterated until convergence or further **niter2** C-steps

- Combining **partially correct** information from random initializations can lead to effective **ensemble initializations** [Alvarez-Esteban et al 2025+]:

- ① Partitions  $\{\mathcal{C}_b\}_{b=1}^{nstart}$  resulting from random initializations
- ② Build an affinity matrix  $A$  ( $n \times n$ ) with:

$$A_{ii'} = \frac{1}{nstart} \# \{ b : x_i \text{ and } x_{i'} \text{ co-clustered (and non-trimmed) in } \mathcal{C}_b \}$$

- ③ Exploit the information in matrix  $A$  for new initializations

# Difficulties in the increasing number of parameters

Robust  
clustering in  
higher  
dimensions

L.A. García-  
Escudero

Robust  
clustering and  
trimming

TCLUST high  
dimensions

RLG method

Trimmed  
HDDC

Other issues

Conclusions

- A **huge number of parameters**:

$$(G - 1) + G \cdot p + G \cdot \frac{(p + 1)p}{2}.$$

- A **huge number of parameters**:

$$(G - 1) + G \cdot p + G \cdot \frac{(p + 1)p}{2}.$$

- ◇ What makes TCLUST attractive and flexible for low  $p$  leads to failure when  $p$  grows...

- A **huge number of parameters**:

$$(G - 1) + G \cdot p + G \cdot \frac{(p + 1)p}{2}.$$

- ◇ What makes TCLUST attractive and flexible for low  $p$  leads to failure when  $p$  grows...
- ◇ Even more problematic when  $n$  is small relative to  $p$

- A **huge number of parameters**:

$$(G - 1) + G \cdot p + G \cdot \frac{(p + 1)p}{2}.$$

- ◇ What makes TCLUST attractive and flexible for low  $p$  leads to failure when  $p$  grows...
- ◇ Even more problematic when  $n$  is small relative to  $p$
- The eigenvalue ratio constraint **“regularizes”** the target function when  $c$  is small. However, small  $c$ 's ( $c \approx 1$ ) enforce **trimmed  $k$ -means-type** results

- A **huge number of parameters**:

$$(G - 1) + G \cdot p + G \cdot \frac{(p + 1)p}{2}.$$

- ◇ What makes TCLUST attractive and flexible for low  $p$  leads to failure when  $p$  grows...
- ◇ Even more problematic when  $n$  is small relative to  $p$
- The eigenvalue ratio constraint **“regularizes”** the target function when  $c$  is small. However, small  $c$ 's ( $c \approx 1$ ) enforce **trimmed  $k$ -means-type** results
- More sophisticated constraints [G-E, Mayo-Iscar and Riani 2020, 2022]

# Robust Linear Grouping method





# Dimension reduction and clustering

- Applying PCA and then clustering (tandem approach) is not the best strategy [Chang 1983]

# Dimension reduction and clustering

- Applying PCA and then clustering (tandem approach) is not the best strategy [Chang 1983]
- Perform **clustering** and **dimension reduction** simultaneously

# Dimension reduction and clustering

- Applying PCA and then clustering (tandem approach) is not the best strategy [Chang 1983]
- Perform **clustering** and **dimension reduction** simultaneously
- Observations clustered around  $G$  **affine subspaces**:
  - ◇ Mixtures of PPCA [Tipping and Bishop 1997]
  - ◇ Mixtures of Factor Analyzers [Ghahramani and Hinton 1997; McLachlan and Peel 2000]

- The **Robust Linear Grouping (RLG)** searches for  $G$  **affine subspaces**  $\mathcal{B}_1, \dots, \mathcal{B}_G$  with intrinsic dimensions  $q_1, \dots, q_G$  and a partition  $\{R_0, R_1, \dots, R_G\}$  with  $\#R_0 = \lfloor n\alpha \rfloor$  minimizing

$$\sum_{g=1}^G \sum_{i \in R_g} \|x_i - \text{Pr}_{\mathcal{B}_g}(x_i)\|^2$$

- ◇  $\text{Pr}_{\mathcal{B}}(x) \rightsquigarrow$  orthogonal projection of  $x \in \mathbb{R}^p$  onto subspace  $\mathcal{B}$

- The **Robust Linear Grouping (RLG)** searches for  $G$  **affine subspaces**  $\mathcal{B}_1, \dots, \mathcal{B}_G$  with intrinsic dimensions  $q_1, \dots, q_G$  and a partition  $\{R_0, R_1, \dots, R_G\}$  with  $\#R_0 = \lfloor n\alpha \rfloor$  minimizing

$$\sum_{g=1}^G \sum_{i \in R_g} \|x_i - \text{Pr}_{\mathcal{B}_g}(x_i)\|^2$$

- ◇  $\text{Pr}_{\mathcal{B}}(x) \rightsquigarrow$  orthogonal projection of  $x \in \mathbb{R}^p$  onto subspace  $\mathcal{B}$
- ◇ LTS-PCA when  $G = 1$  [Maronna 2005; Croux et al 2017]

- The **Robust Linear Grouping (RLG)** searches for  $G$  **affine subspaces**  $\mathcal{B}_1, \dots, \mathcal{B}_G$  with intrinsic dimensions  $q_1, \dots, q_G$  and a partition  $\{R_0, R_1, \dots, R_G\}$  with  $\#R_0 = \lfloor n\alpha \rfloor$  minimizing

$$\sum_{g=1}^G \sum_{i \in R_g} \|x_i - \text{Pr}_{\mathcal{B}_g}(x_i)\|^2$$

- ◇  $\text{Pr}_{\mathcal{B}}(x) \rightsquigarrow$  orthogonal projection of  $x \in \mathbb{R}^p$  onto subspace  $\mathcal{B}$
- ◇ LTS-PCA when  $G = 1$  [Maronna 2005; Croux et al 2017]
- ◇ Trimmed  $k$ -means when  $q_1 = \dots = q_G = 0$

## 1 *Random initializations:*

$\sum_{g=1}^G (q_g + 1)$  random observ.  $\rightsquigarrow$  Affine suspases  $\mathcal{B}_1^0, \dots, \mathcal{B}_G^0$

## 1 *Random initializations:*

$\sum_{g=1}^G (q_g + 1)$  random observ.  $\rightsquigarrow$  Affine suspases  $\mathcal{B}_1^0, \dots, \mathcal{B}_G^0$

## 2 *C-steps:*

### 2.1 Take

$$D_{ig} = \|x_i - \text{Pr}_{\mathcal{B}_g^{t-1}}(x_i)\|^2 \text{ and } D_i = \min_{g=1, \dots, G} D_{ig}$$

and  $D_{(1)} \leq D_{(2)} \leq \dots \leq D_{(n)}$  to get

$$R_g = \{i : D_{ig} = D_i \text{ and } D_i \leq D_{([n(1-\alpha)])}\}$$

### 2.2 Update $\mathcal{B}_g^t$ based on $\mathbf{PCA}_{q_g}\{x_i : i \in R_g\}$



### ① *Random initializations:*

$\sum_{g=1}^G (q_g + 1)$  random observ.  $\rightsquigarrow$  Affine suspases  $\mathcal{B}_1^0, \dots, \mathcal{B}_G^0$

### ② *C-steps:*

#### 2.1 Take

$$D_{ig} = \|x_i - \text{Pr}_{\mathcal{B}_g^{t-1}}(x_i)\|^2 \text{ and } D_i = \min_{g=1, \dots, G} D_{ig}$$

and  $D_{(1)} \leq D_{(2)} \leq \dots \leq D_{(n)}$  to get

$$R_g = \{i : D_{ig} = D_i \text{ and } D_i \leq D_{([n(1-\alpha)])}\}$$

#### 2.2 Update $\mathcal{B}_g^t$ based on $\mathbf{PCA}_{q_g}\{x_i : i \in R_g\}$

### ③ *Output* that with the minimum value of the target function

### 1 *Random initializations:*

$\sum_{g=1}^G (q_g + 1)$  random observ.  $\rightsquigarrow$  Affine suspases  $\mathcal{B}_1^0, \dots, \mathcal{B}_G^0$

### 2 *C-steps:*

#### 2.1 Take

$$D_{ig} = \|x_i - \text{Pr}_{\mathcal{B}_g^{t-1}}(x_i)\|^2 \text{ and } D_i = \min_{g=1, \dots, G} D_{ig}$$

and  $D_{(1)} \leq D_{(2)} \leq \dots \leq D_{(n)}$  to get

$$R_g = \{i : D_{ig} = D_i \text{ and } D_i \leq D_{([n(1-\alpha)])}\}$$

#### 2.2 Update $\mathcal{B}_g^t$ based on $\mathbf{PCA}_{q_g}\{x_i : i \in R_g\}$

### 3 *Output* that with the minimum value of the target function

- It can be applied using `r1g()` function in `tclust` package

- Isotropic orthogonal errors and troubles with intersecting subspaces:

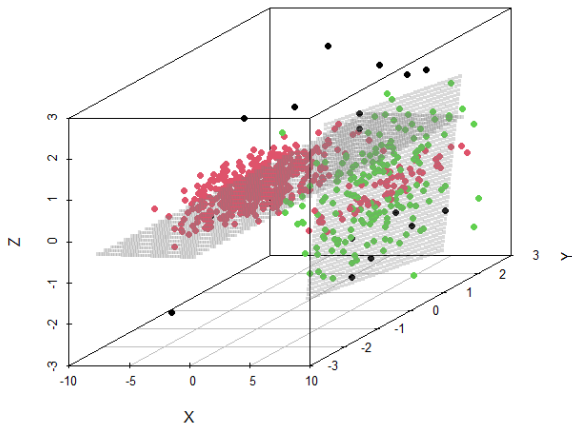


Figure: Troubles with intersecting subspaces ( $q_1 = q_2 = 2$ )

# Trimmed HDDC



- A compromise between TCLUST and RLG

- A **compromise between TCLUST and RLG**
- It arises from the **HDDC** (High Dimensional Data Clustering) approach [Bouveyron, Girard and Schmid 2007; Bouveyron and Brunet-Saumard 2014], where

$$\Sigma_g = U_g \Delta_g U_g',$$

- ◇  $U_g$  is the orthonormal matrix with the eigenvectors of  $\Sigma_g$
- ◇  $\Delta_g$  is a diagonal matrix with the sorted eigenvalues, but with a special **parsimonious** structure on the  $\Delta_g$  matrix.

- $\Delta_g$  has the form

$$\Delta_g = \left( \begin{array}{ccc|ccc} \lambda_{g1} & & & & & \\ & \ddots & & & & \\ & & 0 & & & \\ \hline 0 & & \lambda_{gq_g} & & & \\ & & & \lambda_g & & 0 \\ & & & & \ddots & \\ & & 0 & & & \lambda_g \end{array} \right) \left\{ \begin{array}{l} q_g \\ (p - q_g) \end{array} \right.$$

where:

- ◇  $\lambda_{gl} \geq \lambda_g$  for  $l = 1, \dots, q_g$  and  $g = 1, \dots, G$
- ◇  $q_g < p$  for  $g = 1, \dots, G$

- Additional **eigenvalues-ratio constraints** imposed:

$$\frac{\max_{g=1,\dots,G,j=1,\dots,q_g} \lambda_{gj}}{\min_{g=1,\dots,G,j=1,\dots,q_g} \lambda_{gj}} \leq c_1,$$

and

$$\frac{\max_{g=1,\dots,G} \lambda_g}{\min_{g=1,\dots,G} \lambda_g} \leq c_2$$

- The second constraint is the most relevant one



- Let  $\mathcal{B}_g$  be the affine subspace passing through  $\mu_g$  and spanned by the first  $q$  columns of  $U_g$  denoted as  $u_{g1}, \dots, u_{gq_g}$

- We have

$$\begin{aligned} \log[\pi_g \phi(x_i; m_g, \Sigma_g)] = & \log(\pi_g) - \frac{1}{2} \left( \|\text{Pr}_{\mathcal{B}_g}(x_i) - \mu_g\|_{\mathcal{B}_g}^2 \right. \\ & + \frac{1}{\lambda_g} \|x_i - \text{Pr}_{\mathcal{B}_g}(x_i)\|^2 \\ & + \sum_{l=1}^{q_g} \log(\lambda_{gl}) + (p - q_g) \log(\lambda_g) \\ & \left. + p \log(2\pi) \right) \end{aligned}$$

- Let  $\mathcal{B}_g$  be the affine subspace passing through  $\mu_g$  and spanned by the first  $q$  columns of  $U_g$  denoted as  $u_{g1}, \dots, u_{gq_g}$

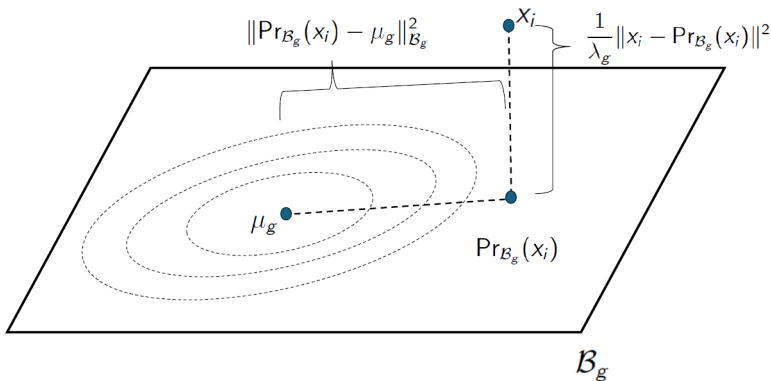
- We have

$$\begin{aligned} \log[\pi_g \phi(x_i; m_g, S_g)] = & \log(\pi_g) - \frac{1}{2} \left( \sum_{l=1}^{q_g} \frac{\langle x_i - \mu_g, u_{gl} \rangle^2}{\lambda_{gl}} \right. \\ & + \frac{1}{\lambda_g} \left\| x_i - \mu_g - \sum_{l=1}^{q_g} \langle x_i - \mu_g, u_{gl} \rangle u_{gl} \right\|^2 \\ & + \sum_{l=1}^{q_g} \log(\lambda_{gl}) + (p - q_g) \log(\lambda_g) \\ & \left. + p \log(2\pi) \right) \rightsquigarrow \text{Only first } q_g \text{ eigenv.!!!} \end{aligned}$$

- Let  $\mathcal{B}_g$  be the affine subspace passing through  $\mu_g$  and spanned by the first  $q$  columns of  $U_g$  denoted as  $u_{g1}, \dots, u_{gq_g}$
- We have

$$\begin{aligned} \log[\pi_g \phi(x_i; m_g, S_g)] = & \log(\pi_g) - \frac{1}{2} \left( \underbrace{\|\text{Pr}_{\mathcal{B}_g}(x_i) - \mu_g\|_{\mathcal{B}_g}^2}_{\text{TCLUST type in } \mathcal{B}_g} \right. \\ & + \frac{1}{\lambda_g} \underbrace{\|x_i - \text{Pr}_{\mathcal{B}_g}(x_i)\|^2}_{\text{RLG type in } \mathcal{B}_g^\perp} \\ & + \sum_{l=1}^{q_g} \log(\lambda_{gl}) + (p - q_g) \log(\lambda_g) \\ & \left. + p \log(2\pi) \right) \end{aligned}$$

## Summary:



**Figure:** Trimming  $x_i$  with large orthogonal residuals or large distances after projection onto the approximating affine subspace

- Trimmed **classification/mixture** likelihoods as

$$\sum_{g=1}^G \sum_{i \in R_g} \log [\pi_g \phi(x_i; \mu_g, \Sigma_g)]$$

and

$$\sum_{i \in R} \log \left[ \sum_{g=1}^G \pi_g \phi(x_i; \mu_g, \Sigma_g) \right]$$

that (in both cases) require (in the M-step) the maximization of **completed likelihoods**:

$$\sum_{i=1}^n \sum_{g=1}^G z_{ig} \sum_{i \in R_g} \log [\pi_g \phi(x_i; \mu_g, \Sigma_g)]$$

① *Random initializations: To be explained later...*

① *Random initializations:* To be explained later...

② *C-steps:*

2.1 *Update  $z_{ig}$ 's:* Let

$$D_{ig} = \pi_g^{t-1} \phi(x_i, \mu_g^{t-1}, \overbrace{\{u_{gl}^{t-1}, \lambda_{gl}^{t-1}, \lambda_g^{t-1}\}_{l=1}^{q_g}}^{\text{Relevant info from } \Sigma_g^{t-1}}),$$

$$D_i = \max_{g=1, \dots, G} D_{ig} \text{ (clasiff.)}, \quad D_i = \sum_{g=1}^G D_{ig} \text{ (mixt.) and}$$

$$z_{ig} = \begin{cases} 1 & \text{if } D_{ig} = D_i \text{ (clasiff.) or } z_{ig} = \frac{D_{ig}}{\sum_{g=1}^G D_{ig}} \text{ (mixt.)}, \\ 0 & \text{otherwise} \end{cases}$$

but  $z_{ig} = 0$ , for every  $g = 1, \dots, G$ , if  $D_i \leq D_{([n\alpha])}$

## 2.2 *Update parameters:*

**2.2.1** Update weights  $\pi_g^t = \frac{\sum_{g=1}^G z_{ig}}{[n(1 - \alpha)]}$

**2.2.2** Update means  $\mu_g^t = \frac{\sum_{g=1}^G z_{ig} x_i}{\sum_{g=1}^G z_{ig}}$

**2.2.3** Compute  $S_g = \frac{1}{\sum_{g=1}^G z_{ig}} \sum_{g=1}^G z_{ig} (x_i - \mu_g^l)(x_i - \mu_g^t)'$  and obtain its eigenvectors  $u_{g1}^t, \dots, u_{gq_g}^t$  associated to the  $q_g$ -th largest eigenvalues  $d_{g1}, \dots, d_{gq_g}$  and obtain

$$d_g = \frac{1}{p - q_g} (\text{trace}(S_g) - \sum_{l=1}^{q_g} d_{gl}).$$



### 2.2.4 Impose constraints on

$$\{\{d_{gl}\}_{l=1,\dots,q_g}\}_{g=1}^G \text{ and } \{d_g\}_{g=1}^G$$

by applying twice the optimal truncation operation described for TCLUST with constants  $c_1$  and  $c_2$  and return

$$\{\{\lambda_{gl}^t\}_{l=1,\dots,q_g}\}_{g=1}^G \text{ and } \{\lambda_g^t\}_{g=1}^G$$

③ *Output* that with the maximum value of the target function

### 2.2.4 Impose constraints on

$$\{\{d_{gl}\}_{l=1,\dots,q_g}\}_{g=1}^G \text{ and } \{d_g\}_{g=1}^G$$

by applying twice the optimal truncation operation described for TCLUST with constants  $c_1$  and  $c_2$  and return

$$\{\{\lambda_{gl}^t\}_{l=1,\dots,q_g}\}_{g=1}^G \text{ and } \{\lambda_g^t\}_{g=1}^G$$

- ③ *Output* that with the maximum value of the target function
- Only eigenvalues/eigenvectors associated to the  $q_g$  largest eigenvalues required (Arnoldi's method...)

### 2.2.4 Impose constraints on

$$\{\{d_{gl}\}_{l=1,\dots,q_g}\}_{g=1}^G \text{ and } \{d_g\}_{g=1}^G$$

by applying twice the optimal truncation operation described for TCLUST with constants  $c_1$  and  $c_2$  and return

$$\{\{\lambda_{gl}^t\}_{l=1,\dots,q_g}\}_{g=1}^G \text{ and } \{\lambda_g^t\}_{g=1}^G$$

③ *Output* that with the maximum value of the target function

- Only eigenvalues/eigenvectors associated to the  $q_g$  largest eigenvalues required (Arnoldi's method...)
- TCLUST as a particular case if  $q_g = p - 1$  and  $c_1 = c_2$  and connections with the trimmed MFA [G-E et al 2016]

- Select a random subsample of size  $\sum_{g=1}^G (q_g + 2) (\ll G \times (p+1))$

- Select a random subsample of size  $\sum_{g=1}^G (q_g + 2) (\ll G \times (p+1))$
- The minimal set of observations to initialize all parameters...

- Select a random subsample of size  $\sum_{g=1}^G (q_g + 2) (\ll G \times (p+1))$
- The minimal set of observations to initialize all parameters...
- From observations  $x_{i_1}, \dots, x_{i_{q_g+2}}$  (in general position):
  - ◊ Take  $\mu_g^0 = \mathbf{avg}\{x_{i_1}, \dots, x_{i_{q_g+2}}\}$
  - ◊ Take  $u_{g1}^0, \dots, u_{gq_g}^0$  and  $d_{g1}, \dots, d_{gq_g}$  associated to the  $q_g$  largest eigenvalues of  $\mathbf{cov}\{x_{i_1}, \dots, x_{i_{q_g+2}}\}$
  - ◊  $d_g$  equal to the smallest eigenvalue divided by  $p - q_g$
  - ◊ If  $\mathfrak{X}$  is the  $(q_g+2) \times p$  matrix whose columns are these observations centred with  $\mu_g^0$  then, instead of using  $p \times p$  matrix  $\mathfrak{X}'\mathfrak{X}$ , use the  $(q_g+2) \times (q_g+2)$  matrix  $\mathfrak{X}\mathfrak{X}'$ ...

- **tHDDC** with  $q_1 = q_2 = q_3 = 8$ ,  $c_1 = 5$ ,  $c_2 = 2$  and  $\alpha = 0.05$  for the **Digits data**:

	0	1	2	3
3	20	29	23	586
5	23	3	510	20
8	44	484	7	7

- **tHDDC** with  $q_1 = q_2 = q_3 = 8$ ,  $c_1 = 5$ ,  $c_2 = 2$  and  $\alpha = 0.05$  for the **Digits data**:

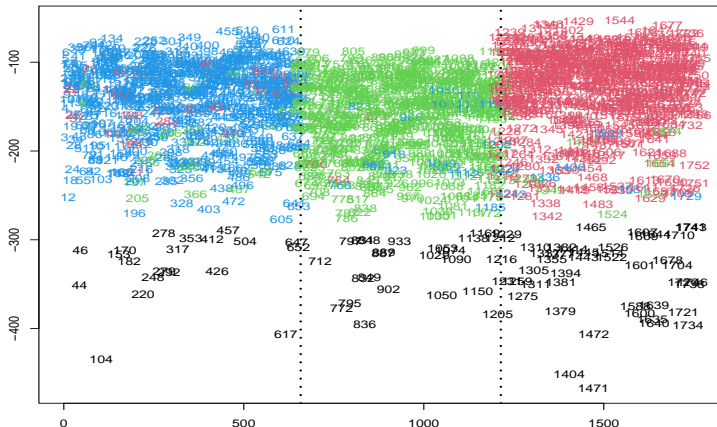
	0	1	2	3
3	20	29	23	586
5	23	3	510	20
8	44	484	7	7

- Better performance than TCLUST and RLG



## Digits data: outliers

- If  $D_{ig} = \pi_g \phi(x_i; \mu_g, \Sigma_g)$  and  $D_i = \max_{g=1, \dots, G} D_{ig}$ , the trimmed observations are those with the smallest  $D_i$ :



**Figure:**  $D_i$  for  $i = 1, \dots, 1756$  with cluster assignments in colors and black for trimmed

- Some of the observations with the smallest  $D_i$  values (trimmed ones):

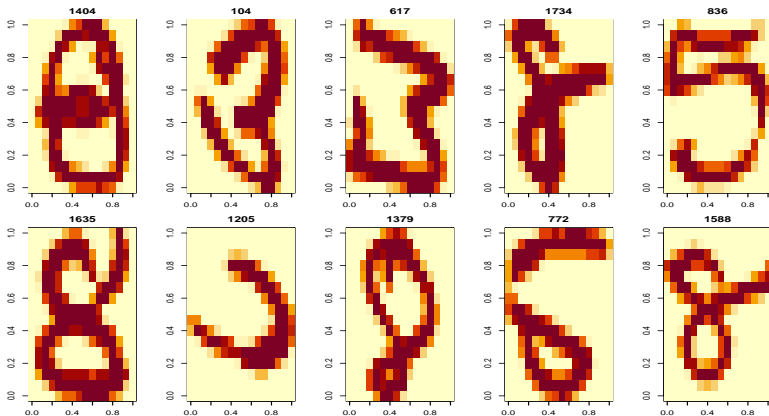


Figure: Some trimmed units

# Digits data: estimated centers

- The estimated location vectors  $\mu_1$ ,  $\mu_2$  and  $\mu_3$ :

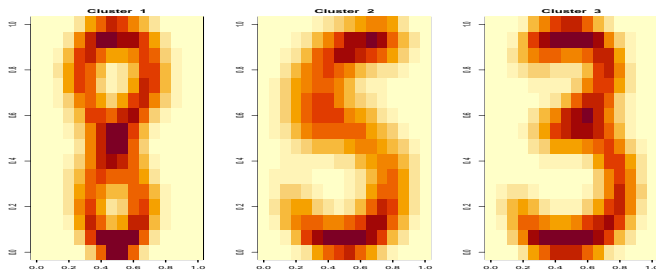
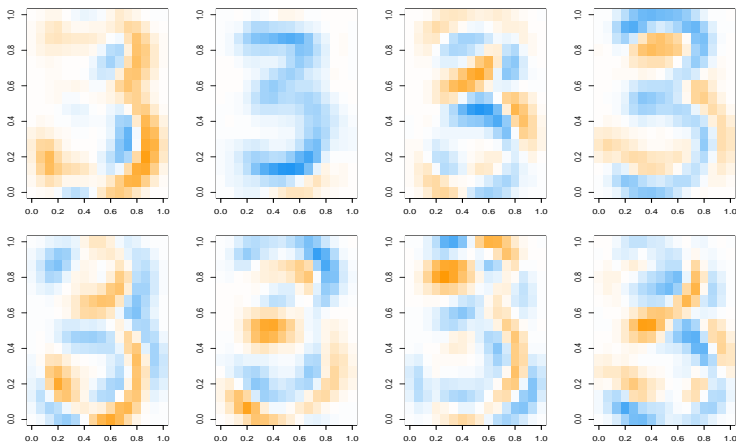


Figure: Estimated clusters' location vectors

- $u_{g1}, u_{g2}, \dots, u_{gq_g}$  loadings interpreted as “variations modes”:



**Figure:** Loadings ( $q_3 = 8$ ) for the cluster including the 3's: ■ for positive weights and ■ for negative ones



- **Score distances** as

$$SD_{ig} = \|\text{Pr}_{\mathcal{B}_g}(x_i) - \mu_g\|_{\mathcal{B}_g} = \sqrt{\sum_{l=1}^{q_g} \frac{t_{ig,l}^2}{\lambda_{g,l}}}$$

- **Score distances** as

$$SD_{ig} = \|\text{Pr}_{\mathcal{B}_g}(x_i) - \mu_g\|_{\mathcal{B}_g} = \sqrt{\sum_{l=1}^{q_g} \frac{t_{ig,l}^2}{\lambda_{g,l}}}$$

- **Orthogonal distances** as

$$OD_{ig} = \frac{1}{\sqrt{\lambda_g}} \|x_i - \text{Pr}_{\mathcal{B}_g}(x_i)\| = \frac{1}{\sqrt{\lambda_g}} \left\| x_i - \mu_g - \sum_{l=1}^{q_g} t_{ig,l} u_{g,l} \right\|$$

- **Score distances** as

$$SD_{ig} = \|\text{Pr}_{\mathcal{B}_g}(x_i) - \mu_g\|_{\mathcal{B}_g} = \sqrt{\sum_{l=1}^{q_g} \frac{t_{ig,l}^2}{\lambda_{g,l}}}$$

- **Orthogonal distances** as

$$OD_{ig} = \frac{1}{\sqrt{\lambda_g}} \|x_i - \text{Pr}_{\mathcal{B}_g}(x_i)\| = \frac{1}{\sqrt{\lambda_g}} \left\| x_i - \mu_g - \sum_{l=1}^{q_g} t_{ig,l} u_{g,l} \right\|$$

- Computed for  $A_g = \{i : D_{ig} = D_i\}$  (trimmed observations also included) and **cutoffs**  $\sqrt{\chi_{q_g;0.975}^2}$  for the  $SD_{ig}$  and  $(\hat{\mu}_{OD} + \hat{\sigma}_{OD} z_{0.975})^{3/2}$  where  $\hat{\mu}_{OD}$  and  $\hat{\sigma}_{OD}^2$  are the robust mean and variance based on MCD on  $\{OD_{ig}^{2/3}\}$  for the  $OD_{ig}$  [Hubert et al 2005]



## Digits data: diagnostic plot example

- Score and orthogonal distances:

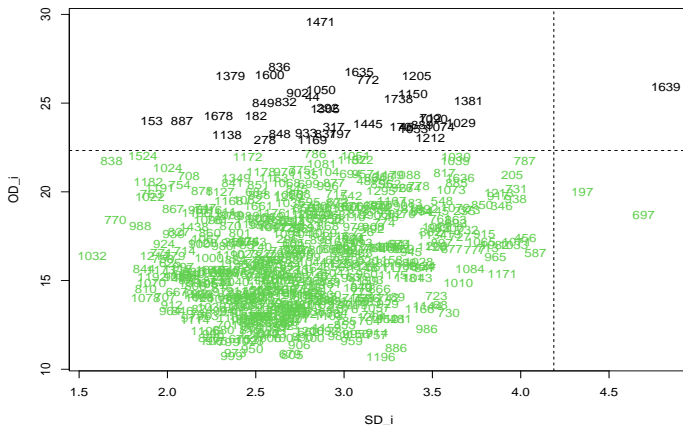


Figure:  $SD_{ig}$  and  $OD_{ig}$  for  $g$  corresponding to the cluster with the 3's

- Some of the atypical digits detected in the previous plot:

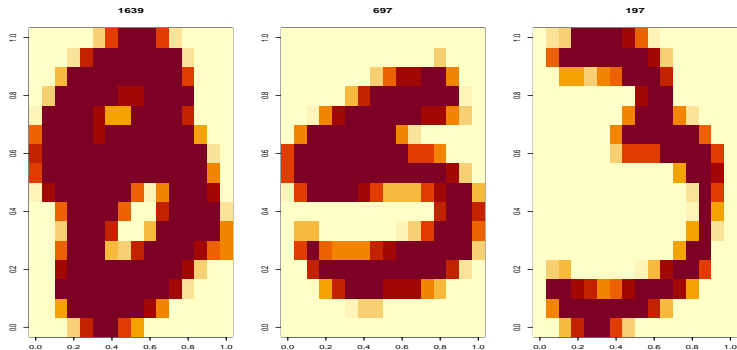


Figure: Some highlighted outliers using  $SD_{ig}$  and  $OD_{ig}$  for the cluster with the 3's

- From the  $D_{ig} = \pi_g \phi(x_i; \mu_g, \Sigma_g)$ , let

$$D_{i(1)} \leq D_{i(2)} \leq \dots \leq D_{i(G)}$$

and define **Discriminant Factors** as

$$DF(i) = \log \frac{D_{i(G-1)}}{D_{i(G)}} \text{ if } x_i \text{ is not trimmed}$$

- From the  $D_{ig} = \pi_g \phi(x_i; \mu_g, \Sigma_g)$ , let

$$D_{i(1)} \leq D_{i(2)} \leq \dots \leq D_{i(G)}$$

and define **Discriminant Factors** as

$$\text{DF}(i) = \log \frac{D_{i(G-1)}}{D_{i(G)}} \text{ if } x_i \text{ is not trimmed}$$

- If  $D^* = D_{([n\alpha])}$  is the cutoff to label outliers then

$$\text{DF}(i) = \log \frac{D_i}{D^*} \text{ if } x_i \text{ trimmed}$$

- From the  $D_{ig} = \pi_g \phi(x_i; \mu_g, \Sigma_g)$ , let

$$D_{i(1)} \leq D_{i(2)} \leq \dots \leq D_{i(G)}$$

and define **Discriminant Factors** as

$$\text{DF}(i) = \log \frac{D_{i(G-1)}}{D_{i(G)}} \text{ if } x_i \text{ is not trimmed}$$

- If  $D^* = D_{([n\alpha])}$  is the cutoff to label outliers then

$$\text{DF}(i) = \log \frac{D_i}{D^*} \text{ if } x_i \text{ trimmed}$$

- $\text{DF}(i) \leq 0$  but  $\text{DF}(i) \simeq 0$  for the most **doubtful** assignment decisions.

- Silhouette plot:

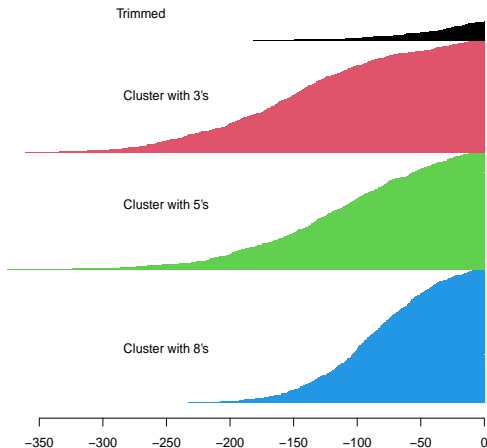


Figure: Silhouette plot with the  $DF(i)$  values

- $q_1, \dots, q_G$  (even  $G \dots$ ),  $\alpha$ ,  $c_1$  and  $c_2$ ???

# Choice of hyperparameters

- $q_1, \dots, q_G$  (even  $G \dots$ ),  $\alpha$ ,  $c_1$  and  $c_2$ ???
- A **complex** problem... as in others (robust) clustering methods



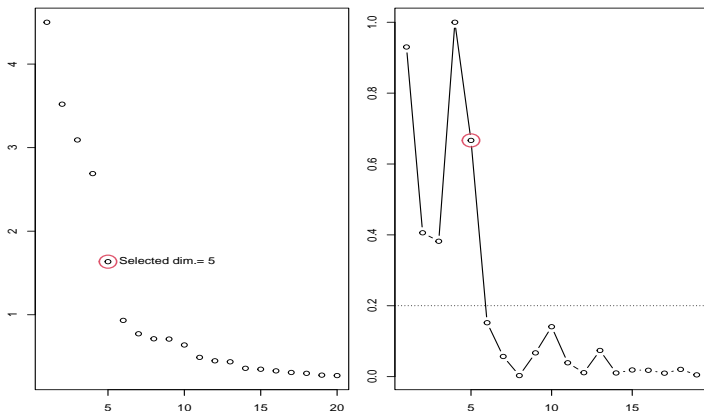
# Choice of hyperparameters

- $q_1, \dots, q_G$  (even  $G \dots$ ),  $\alpha$ ,  $c_1$  and  $c_2$ ???
- A **complex** problem... as in others (robust) clustering methods
- Cannot be fully automated because user decisions are often required

- $q_1, \dots, q_G$  (even  $G \dots$ ),  $\alpha$ ,  $c_1$  and  $c_2$ ???
- A **complex** problem... as in others (robust) clustering methods
- Cannot be fully automated because user decisions are often required
- The intrinsic dimensions  $q_1, \dots, q_G$  can be treated as **inner parameters to be estimated** within the iterative steps of the algorithm

Catell procedure estimating  $q_g$ 's

- **Catell procedure** based on differences  $\lambda_{g,l+1} - \lambda_{g,l}$  of the sorted eigenvalues of  $S_g$  (with an upper bound  $q_{\max}$ ):



**Figure:** Estimate  $q_g$  as largest index where (normalized) differences exceed a threshold  $\text{tresh}$ :  $q_g = 5$  selected with  $\text{tresh}=0.2$

- Estimating the  $q_g$ 's requires a **BIC-type** [Cerioli et al 2018] modified **target function**:

- Estimating the  $q_g$ 's requires a **BIC-type** [Cerioli et al 2018] modified **target function**:

$$-2 \times \text{trimmed log-likelihood}(q_1, \dots, q_G) + \text{complexity penalty}$$

Catell procedure estimating  $q_g$ 's

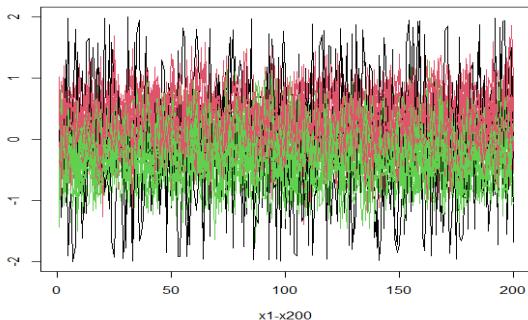
- Estimating the  $q_g$ 's requires a **BIC-type** [Cerioli et al 2018] modified **target function**:

$-2 \times \text{trimmed log-likelihood}(q_1, \dots, q_G) + \text{complexity penalty}$   
with **complexity penalty** equal to

$$\log([n(1-\alpha)]) \left[ \underbrace{G-1}_{\text{weights}} + \underbrace{Gp}_{\text{means}} + 1 + \underbrace{\left( \sum_{g=1}^G q_g - 1 \right) \left( 1 - \frac{1}{c_1} \right)}_{\text{the largest eigenvalues}} \right. \\ \left. + \underbrace{1 + (G-1) \left( 1 - \frac{1}{c_2} \right)}_{\text{the smallest ones}} + \underbrace{\sum_{g=1}^G \left( q_g p - \frac{q_g(q_g-1)}{2} \right)}_{\text{orthonormal eigenvec.}} \right]$$

# Example I: estimating $q_g$ 's

- **Simulated data** set with  $n = 1000$  and  $p = 200$  where
  - ◇  $\mu_1 = 0.2 \cdot \mathbf{1}_{200}$  and  $\Sigma_1 = \text{diag}(5, 4, 3, 2, 1, 0.1, \dots, 0.1)$  : 57% rows
  - ◇  $\mu_2 = -0.2 \cdot \mathbf{1}_{200}$  and  $\Sigma_2 = \text{diag}(4, 3, 2, 0.15, \dots, 0.15)$  : 38% rows
  - ◇  $x_{ij} \sim \mathcal{U}[-2, 2]$ : 5% rows



**Figure:** Line plots of the simulated dataset with  $q_1 = 5$  and  $q_2 = 3$

Example I: estimating  $q_g$ 's

- **tHDDC** with  $G = 2$ ,  $q_{\max} = 20$ ,  $c_1 = 10$ ,  $c_2 = 2$  and  $\alpha = 0.05$  returns  $\hat{q}_1 = 5$  and  $\hat{q}_2 = 3$ :

$\lambda_{1/}$ :	4.7656	4.1874	3.0381	1.9648	1.0376	$\lambda_1 : 0.0983$
$\lambda_{2/}$ :	3.6971	2.9656	2.0139			$\lambda_2 : 0.1482$

	1	2	0
0	0	0	50
1	570	0	0
2	0	380	0



Example I: estimating  $q_g$ 's

- **tHDDC** with  $G = 2$ ,  $q_{\max} = 20$ ,  $c_1 = 10$ ,  $c_2 = 2$  and  $\alpha = 0.05$  returns  $\hat{q}_1 = 5$  and  $\hat{q}_2 = 3$ :

$\lambda_{1/}$ :	4.7656	4.1874	3.0381	1.9648	1.0376	$\lambda_1 : 0.0983$
$\lambda_{2/}$ :	3.6971	2.9656	2.0139			$\lambda_2 : 0.1482$

	1	2	0
0	0	0	50
1	570	0	0
2	0	380	0

- **TCLUST** with  $G = 2$ ,  $c = 12$  and  $\alpha = 0.05$ :

	1	2	0	
0	0	2	48	
1	502	90	0	# ×10 computing time
2	68	288	2	

## Example II: Digits data

- tHDDC **estimating  $q_g$ 's** with Catell for the **Digits data** returns  $q_1 = q_2 = 10$  and  $q_3 = 9$  (with  $q_{\max} = 20$  and  $\text{tresh}=0.05$ ):

	0	1	2	3
3	25	16	24	593
5	22	8	522	4
8	40	491	2	9

## Example II: Digits data

- tHDDC **estimating  $q_g$ 's** with Catell for the **Digits data** returns  $q_1 = q_2 = 10$  and  $q_3 = 9$  (with  $q_{\max} = 20$  and  $\text{tresh}=0.05$ ):

	0	1	2	3
3	25	16	24	593
5	22	8	522	4
8	40	491	2	9

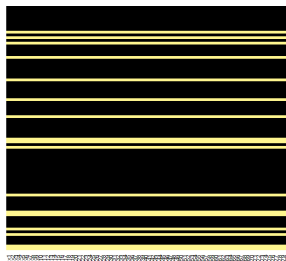
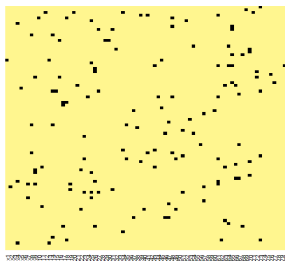
- Even better misclassification rate...

# Other issues



- Robustness against **cellwise contamination** [Alqallaf et al 2009]
  - ◇ **Cases**  $\rightarrow x_i = (x_{i1}, \dots, x_{ip})' \in \mathbb{R}^p$  and **Cells**  $\rightarrow x_{ij} \in \mathbb{R}$

- Robustness against **cellwise contamination** [Alqallaf et al 2009]
  - ◇ **Cases**  $\rightarrow x_i = (x_{i1}, \dots, x_{ip})' \in \mathbb{R}^p$  and **Cells**  $\rightarrow x_{ij} \in \mathbb{R}$
- A lot of useful information sacrificed by casewise trimming:



**Figure:**  $n = 100 \times p = 80$  data-matrix (left) with a 2% of outlying cells and trimmed  $x_i$  with casewise trimming in black (right)

- **Cellwise trimming-TCLUST** [Zaccaria, G-E, Greselin and Mayo-Isacar 2025+]

- **Cellwise trimming-TCLUST** [Zaccaria, G-E, Greselin and Mayo-Isacar 2025+]
  - ◇ Extension to mixture modelling of the cellMCD [Raymaekers and Rousseuw 2024]



- **Cellwise trimming-TCLUST** [Zaccaria, G-E, Greselin and Mayo-Isacar 2025+]
  - ◇ Extension to mixture modelling of the cellMCD [Raymaekers and Rousseuw 2024]
  - ◇ Alternating steps: Detection of outlying cells  $\leftrightarrow$  EM for Gaussian mixtures with missing cells

- **Cellwise trimming-TCLUST** [Zaccaria, G-E, Greselin and Mayo-Isicar 2025+]
  - ◇ Extension to mixture modelling of the cellMCD [Raymaekers and Rousseuw 2024]
  - ◇ Alternating steps: Detection of outlying cells  $\leftrightarrow$  EM for Gaussian mixtures with missing cells
- **Cellwise trimming-RLG** [G-E, Rivera, Mayo-Isicar and Ortega 2021]

- **Cellwise trimming-TCLUST** [Zaccaria, G-E, Greselin and Mayo-Iscar 2025+]
  - ◇ Extension to mixture modelling of the cellMCD [Raymaekers and Rousseuw 2024]
  - ◇ Alternating steps: Detection of outlying cells  $\leftrightarrow$  EM for Gaussian mixtures with missing cells
- **Cellwise trimming-RLG** [G-E, Rivera, Mayo-Iscar and Ortega 2021]
  - ◇ Extension of the cellwise-robust PCA [Maronna and Yohai 2012]

- **Cellwise trimming-TCLUST** [Zaccaria, G-E, Greselin and Mayo-Isicar 2025+]
  - ◇ Extension to mixture modelling of the cellMCD [Raymaekers and Rousseuw 2024]
  - ◇ Alternating steps: Detection of outlying cells  $\leftrightarrow$  EM for Gaussian mixtures with missing cells
- **Cellwise trimming-RLG** [G-E, Rivera, Mayo-Isicar and Ortega 2021]
  - ◇ Extension of the cellwise-robust PCA [Maronna and Yohai 2012]
  - ◇ Alternating Weighted Least Squares (AWLS)

- 1 Given  $\{R_0, R_1, \dots, R_G\}$ , apply AWLS or MacroPCA [Hubert, Rousseeuw and van den Bossche 2019] to update affine subspaces  $\mathcal{B}_g$  and eigenvalues for  $\{x_i : i \in R_g\}$

- 1 Given  $\{R_0, R_1, \dots, R_G\}$ , apply AWLS or MacroPCA [Hubert, Rousseeuw and van den Bossche 2019] to update affine subspaces  $\mathcal{B}_g$  and eigenvalues for  $\{x_i : i \in R_g\}$
- 2  $\Pr_{\mathcal{B}_g}(x_i)$  is needed to update  $\{R_0, R_1, \dots, R_G\} \rightsquigarrow$  Not straightforward for  $\tilde{g} \neq g$  if  $i \in R_g$  (unknown positions of “unreliable” cells...)

- ① Given  $\{R_0, R_1, \dots, R_G\}$ , apply AWLS or MacroPCA [Hubert, Rousseeuw and van den Bossche 2019] to update affine subspaces  $\mathcal{B}_g$  and eigenvalues for  $\{x_i : i \in R_g\}$
- ②  $\Pr_{\mathcal{B}_g}(x_i)$  is needed to update  $\{R_0, R_1, \dots, R_G\} \rightsquigarrow$  Not straightforward for  $\tilde{g} \neq g$  if  $i \in R_g$  (unknown positions of “unreliable” cells...)
- ③ LTS-based predictions for  $\Pr_{\mathcal{B}_{\tilde{g}}}(x_i)$  [G-E, Rivera, Mayo-Isacar and Ortega 2021]

- ① Given  $\{R_0, R_1, \dots, R_G\}$ , apply AWLS or MacroPCA [Hubert, Rousseeuw and van den Bossche 2019] to update affine subspaces  $\mathcal{B}_g$  and eigenvalues for  $\{x_i : i \in R_g\}$
  - ②  $\Pr_{\mathcal{B}_g}(x_i)$  is needed to update  $\{R_0, R_1, \dots, R_G\} \rightsquigarrow$  Not straightforward for  $\tilde{g} \neq g$  if  $i \in R_g$  (unknown positions of “unreliable” cells...)
  - ③ LTS-based predictions for  $\Pr_{\mathcal{B}_{\tilde{g}}}(x_i)$  [G-E, Rivera, Mayo-Isacar and Ortega 2021]
- Computationally intensive approach



- Address **nonlinearity** through **kernelized** versions [Bouveyron, Fauvel and Girard 2013]

- Address **nonlinearity** through **kernelized** versions [Bouveyron, Fauvel and Girard 2013]
- Robust **functional clustering** through functional subspaces [Bouveyron and Jacques 2011]

- Address **nonlinearity** through **kernelized** versions [Bouveyron, Fauvel and Girard 2013]
- Robust **functional clustering** through functional subspaces [Bouveyron and Jacques 2011]
- We just focus on **moderately** (recall the title of the talk!!) high dimensional cases

- Address **nonlinearity** through **kernelized** versions [Bouveyron, Fauvel and Girard 2013]
- Robust **functional clustering** through functional subspaces [Bouveyron and Jacques 2011]
- We just focus on **moderately** (**recall the title of the talk!!**) high dimensional cases
  - ◇ **Noise variables** that do not provide useful information about the clustering structure

- Address **nonlinearity** through **kernelized** versions [Bouveyron, Fauvel and Girard 2013]
- Robust **functional clustering** through functional subspaces [Bouveyron and Jacques 2011]
- We just focus on **moderately** (**recall the title of the talk!!**) high dimensional cases
  - ◇ **Noise variables** that do not provide useful information about the clustering structure
  - ◇ **Variable selection** in robust clustering [Ritter 2014] and **sparsity-based** approaches [Kondo, Salibian-Barrera and Zamar 2016; Brodinova et al 2019; Raymaekers and Zamar 2022]

# Conclusions



- 1 Robust clustering in **high-dimensional** data is of interest

# Conclusions and further directions

- 1 Robust clustering in **high-dimensional** data is of interest
- 2 TCLUST faces limitations in high-dimensional settings



# Conclusions and further directions

- ① Robust clustering in **high-dimensional** data is of interest
- ② TCLUST faces limitations in high-dimensional settings
- ③ Proposed method **tHDDC**, as a compromise between TCLUST and RLG, with initial promising results

# Conclusions and further directions

- 1 Robust clustering in **high-dimensional** data is of interest
- 2 TCLUST faces limitations in high-dimensional settings
- 3 Proposed method **tHDDC**, as a compromise between TCLUST and RLG, with initial promising results
- 4 Plenty of room for further research

Robust  
clustering in  
higher  
dimensions

L.A. García-  
Escudero

Robust  
clustering and  
trimming

TCLUST high  
dimensions

RLG method

Trimmed  
HDDC

Other issues

Conclusions



---

**Universidad de Valladolid**

Thanks for your  
attention!!!!