

# Distance-based probabilistic clustering for functional data

---

*Group members:* Giulia Caruso, Alessio Facincani, Giulia Romani, Pietro Spina, Matteo Vescovi

*Tutors:* Mario Beraha, Riccardo Corradin

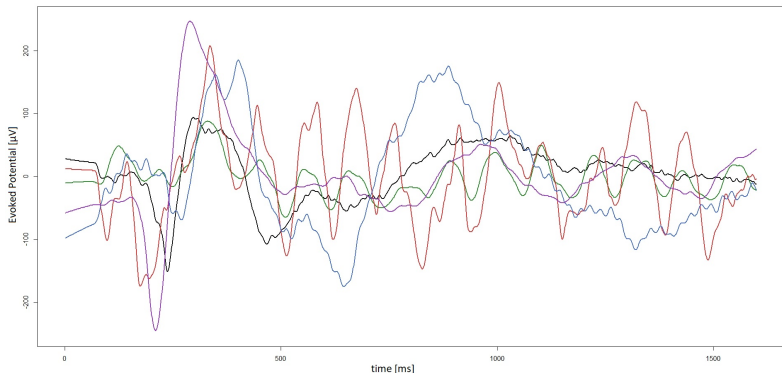


**POLITECNICO**  
MILANO 1863

# Framework of the project

26 **multivariate functional observations**, each component is a function observed at 1600 time points.

RESEARCH GOAL: **Cluster** observations in different groups, first considering only one functional component.



# Model

---

# Prior and posterior

- *Gibbs posterior*:

$$\pi(\mathbf{c}|\lambda, \mathbf{X}) \propto \pi(\mathbf{c}) \exp \{-\lambda \ell(\mathbf{c}; \mathbf{X})\}$$

- *Loss function* under a GB-PPM (Generalized Bayes Product Partition Model) with *Mahalanobis distance*  $M_\alpha(\mathbf{x}_i; \mathbf{X}_k) \geq 0$ :

$$\ell(\mathbf{c}; \mathbf{X}) = \sum_{k=1}^K \sum_{i \in C_k} \mathcal{D}(\mathbf{x}_i; \mathbf{X}_k) = \sum_{k=1}^K \sum_{i \in C_k} M_\alpha(\mathbf{x}_i; \mathbf{X}_k)$$

- *Generalized Bayes posterior under a GB-PPM*:

$$\pi(\mathbf{c}|\lambda, \mathbf{X}) \propto \prod_{k=1}^K \exp \left\{ -\lambda \sum_{i \in C_k} M_\alpha(\mathbf{x}_i; \mathbf{X}_k) \right\}$$

- *Uniform prior (Stirling number of the second kind)*:

$$\pi(\mathbf{c}) = \frac{1}{\mathcal{S}(n, K)}$$

# Mahalanobis distance in a functional context

Given a stochastic process  $X(t) \in \mathbb{L}^2[0, 1]$ ,  $t \in [0, 1]$ :

- *Covariance function  $K$  and operator  $\mathcal{K}$ :*

$$K(s, t) = \text{Cov}(X(s), X(t)) \quad \mathcal{K}f(t) = \int_0^1 K(t, s)f(s)ds$$

- *$\alpha$ -Mahalanobis distance with smoothing parameter  $\alpha > 0$ :*  
( $\forall x, y \in \mathbb{L}^2[0, 1]$ )

$$M_\alpha(x, y)^2 = \|x_\alpha - y_\alpha\|_K^2 = \sum_{j=1}^{\infty} \frac{\lambda_j}{(\lambda_j + \alpha)^2} \langle x - y, e_j \rangle^2$$

- *$\alpha$ -approximation of a function  $x \in \mathbb{L}^2[0, 1]$ :*

$$x_\alpha = \arg \min_{f \in \mathcal{H}(K)} \|x - f\|^2 + \alpha \|f\|_K^2 = \sum_{j=1}^{\infty} \frac{\lambda_j}{(\lambda_j + \alpha)} \langle x, e_j \rangle e_j$$

# Maximum a posteriori estimation (MAP)

Target of inference: optimal and unknown partition  $c_{opt}$

$$\begin{aligned} c_{opt} &= \arg \max_{c : |C|=K} \pi(\mathbf{c} | \lambda, \mathbf{X}) \\ &= \arg \min_{c : |C|=K} \ell(\mathbf{c}; \mathbf{X}) \\ &= \arg \min_{c : |C|=K} \sum_{k=1}^K \sum_{i \in C_k} M_{\alpha}(\mathbf{x}_i; \mathbf{X}_k) \end{aligned}$$

Where the second equality is justified by the uniform prior assumption.

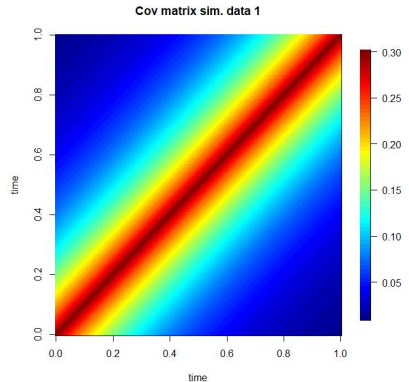
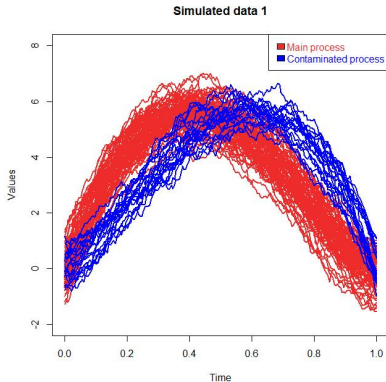
## Posterior inference

---

# Simulated data 1

- **Main process:**  $X(t) = 30t(1 - t)^{3/2} + \epsilon(t)$
- **Contaminated process:**  $X(t) = 30t^{3/2}(1 - t) + \epsilon(t) \quad t \in [0, 1]$

$\epsilon(t) \sim \text{GP}(0, \mathcal{C})$  where  $\mathcal{C}(s, t) = 0.3 \cdot \exp(-|s - t|/0.3)$

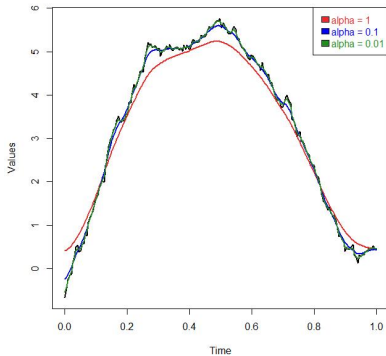




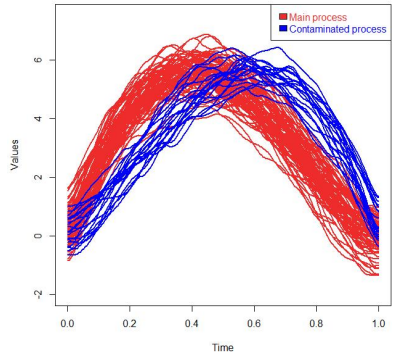
# $\alpha$ -Mahalanobis approximation

*Simulated data 1: best  $\alpha = 10^{-1}$*

Comparison of alpha for sim. data 1: best alpha=0.1



Smoothed sim.data 1



# Algorithm with "fixed" covariance

---

## Algorithm 0: Mahalanobis dist. clustering with fixed covariance

---

**Input** :  $n^\circ$  clusters, covariance matrix,  $\alpha$ , toll, data

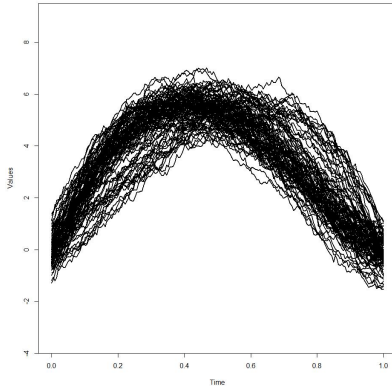
**Output**: Optimal partition (labels, centroids and loss)

```
1 Sample random observations as initial centroids  $\mathbf{m}_1, \dots, \mathbf{m}_K$ 
2 Compute eigenvalues and eigenvectors of the covariance matrix
3 Initialize  $\ell_1(\mathbf{c}; \mathbf{X})$  and  $\ell_2(\mathbf{c}; \mathbf{X})$ 
4 while  $(\ell_1(\mathbf{c}; \mathbf{X}) - \ell_2(\mathbf{c}; \mathbf{X})) > \text{toll}$  do
5      $\ell_1(\mathbf{c}; \mathbf{X}) = \ell_2(\mathbf{c}; \mathbf{X})$ 
6     for  $i=1, \dots, n$  do
7         | Set the cluster indicator  $c_i$  equal to  $k$ , so that  $M_\alpha(\mathbf{x}_i, \mathbf{m}_k)$  is minimum
8     end
9     for  $k=1, \dots, K$  do
10        | Set  $\mathbf{m}_k$  as the functional mean of the observations belonging to cluster  $k$ 
11    end
12    Update  $\ell_2(\mathbf{c}; \mathbf{X})$ 
13 end
```

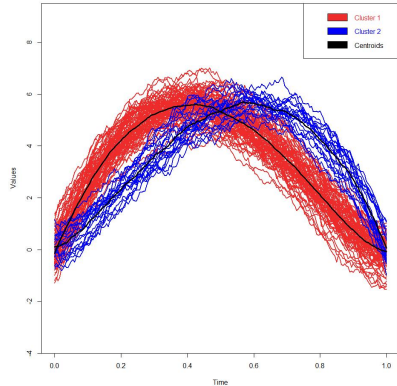
---

# Application on simulated data 1

Main and contaminated processes



Clustered data



# Algorithm with "updated" covariance

---

## Algorithm 1: Mahalanobis dist. clustering with covariance updating

---

**Input** :  $n^\circ$  clusters, covariance matrix,  $\alpha$ , toll, data

**Output**: Optimal partition (labels, centroids and loss)

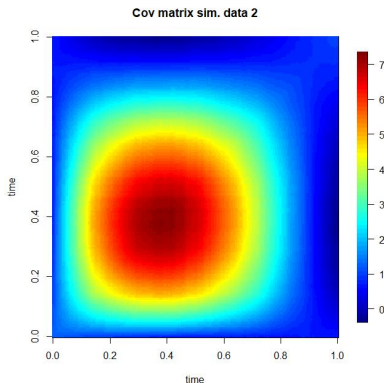
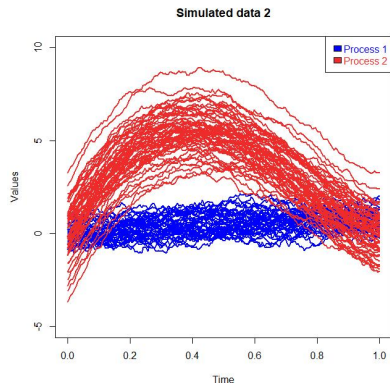
```
1 Sample random observations as initial centroids  $\mathbf{m}_1, \dots, \mathbf{m}_K$ ;  
2 Compute eigenvalues and eigenvectors of the covariance matrix;  
3 Initialize  $\ell_1(\mathbf{c}; \mathbf{X})$  and  $\ell_2(\mathbf{c}; \mathbf{X})$ ;  
4 while  $(\ell_1(\mathbf{c}; \mathbf{X}) - \ell_2(\mathbf{c}; \mathbf{X})) > \text{toll}$  do  
5    $\ell_1(\mathbf{c}; \mathbf{X}) = \ell_2(\mathbf{c}; \mathbf{X})$  ;  
6   for  $i=1, \dots, n$  do  
7     Set the cluster indicator  $c_i$  equal to  $k$ , so that  $M_\alpha(\mathbf{x}_i, \mathbf{m}_k)$  is minimum  
8   end  
9   for  $k=1, \dots, K$  do  
10    Set  $\mathbf{m}_k$  as the functional mean of the observations belonging to cluster  $k$ ;  
11    Set  $\text{cov}_k$  as the covariance matrix of the  $k$ -th cluster and compute its eigenvalues  
    and eigenvectors;  
12  end  
13  Update  $\ell_2(\mathbf{c}; \mathbf{X})$   
14 end
```

---

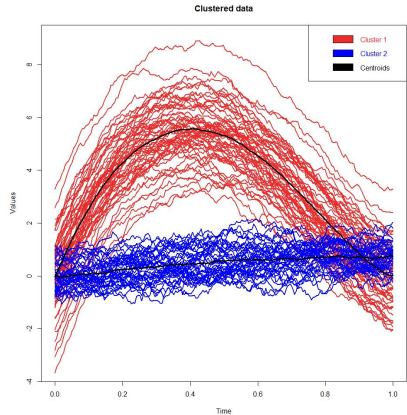
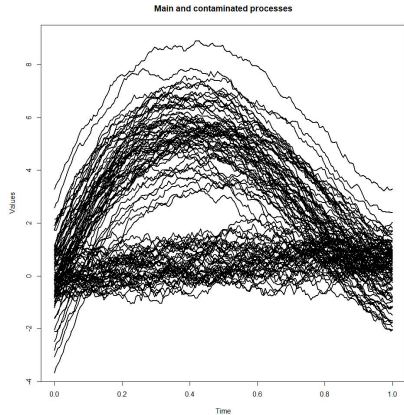
## Simulated data 2

Data clustered into two groups:

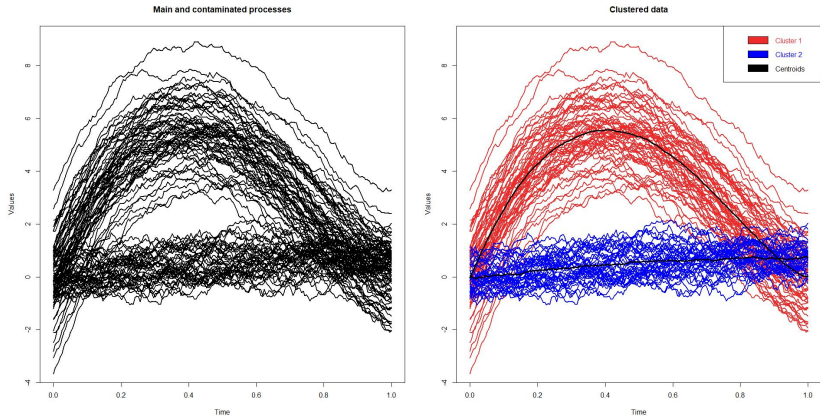
- **First process:**  $X(t) = \sin(t) + \epsilon_1(t)$      $\epsilon_1(t) \sim \text{GP}(0, \mathcal{C}_1)$   
 $\mathcal{C}_1(s, t) = 0.3 \cdot \exp(-|s - t|/0.3)$
- **Second process:**  $X(t) = 30t(1 - t)^{3/2} + \epsilon_2(t)$      $\epsilon_2(t) \sim \text{GP}(0, \mathcal{C}_2)$   
 $\mathcal{C}_2(s, t) = 1.5 \cdot \exp(-|s - t|/3)$     ( $t \in [0, 1]$ )



# Application on simulated data 2



# Application on simulated data 2



**Problem:** the number of clusters has to be given a priori

# Union of similar clusters

---

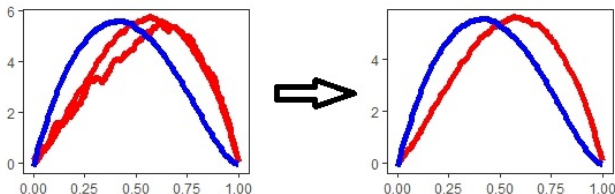
**Algorithm 2:** Union of similar clusters

---

**Input** : optimal centroids  $\mathbf{m}_1, \dots, \mathbf{m}_k$

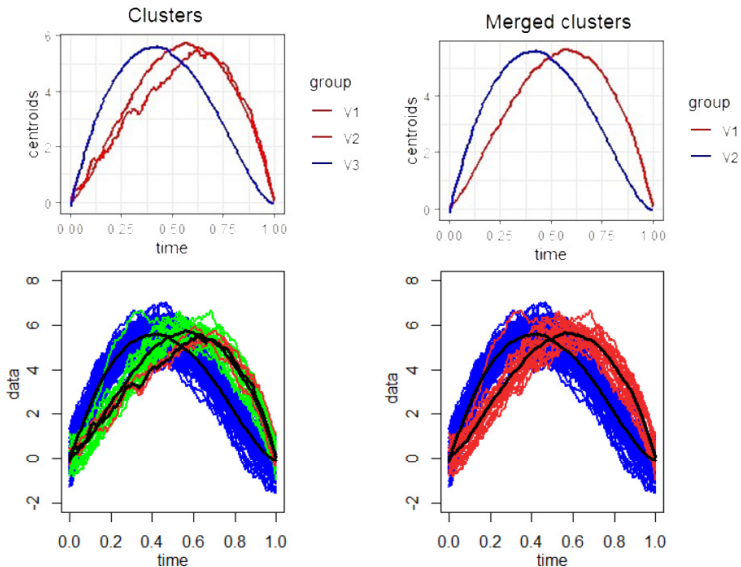
**Output:** united centroids where needed

- 1 Compute distances matrix  $d$  :  $d_{ij} = \|\mathbf{m}_i - \mathbf{m}_j\|_2^2 \quad \forall i, j$ ;
  - 2 Set  $\varepsilon = 0.5 \cdot (\text{median}(d) + \text{mean}(d))$ ;
  - 3 **if**  $d_{ij} < \varepsilon$  **then**
  - 4     | Merge clusters  $i, j$
  - 5 **end**
  - 6 Recompute the centroids and repeat until  $d_{i_{\text{new}}j_{\text{new}}} > \varepsilon_{\text{new}} \quad \forall i_{\text{new}}, j_{\text{new}}$
- 

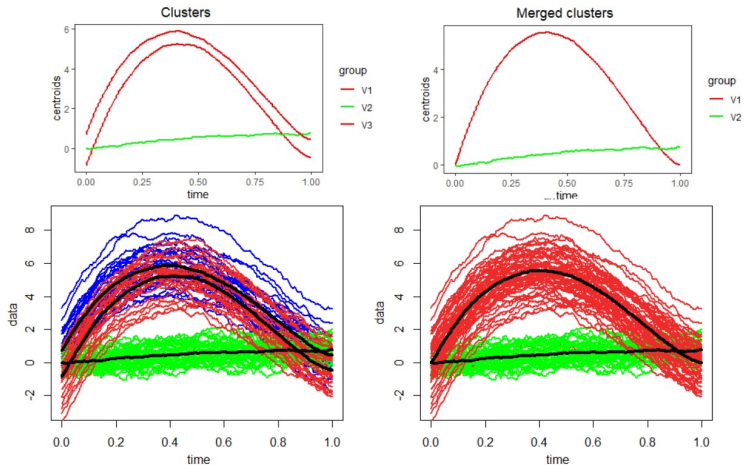




# Union of similar clusters: Simulated data 1



## Union of similar clusters: Simulated data 2



## Next steps

---

- Relax the uniform prior distribution  $\pi(\mathbf{c})$  to a prior that allows "small cluster" penalization
- Application on real data:
  - ◊ Tuning of the smoothing parameter  $\alpha$
- Further development: Gibbs sampling strategy for uncertainty quantification

# References

---

- [1] José R. Berrendero, Beatriz Bueno-Larraz, and Antonio Cuevas. “On Mahalanobis Distance in Functional Settings”. In: *Journal of Machine Learning Research* 21 (2020), pp. 1–33.
- [2] Tommaso Rigon, Amy H. Herring, and David B. Dunson. “A generalized Bayes framework for probabilistic clustering”. In: *arXiv:2006.05451* (2020).
- [3] Sara Wade and Zoubin Ghahramani. “Bayesian Cluster Analysis: Point Estimation and Credible Balls (with Discussion)”. In: *Bayesian Anal.* 13 (2) (2018), pp. 559–626.