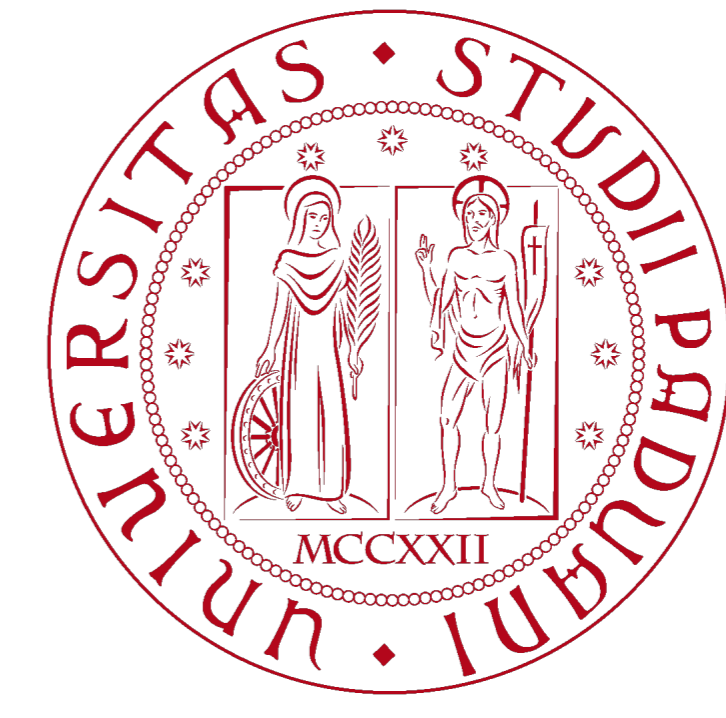


# Local Bayesian clustering for functional data

Giovanni Toto, Antonio Canale

Department of Statistical Sciences, University of Padova



UNIVERSITÀ  
DEGLI STUDI  
DI PADOVA

## Introduction

For standard Euclidean data, clustering is quite straightforward since a point may belong to only a cluster. On the other hand, for functional data the problem becomes more complex since one may perform

- **clustering at global level:** a whole function is assigned to a single cluster,
- **clustering at local level:** the same function may belong to different clusters depending on the point of its domain at which it is evaluated.

**Goal:** performing local clustering for functional data.

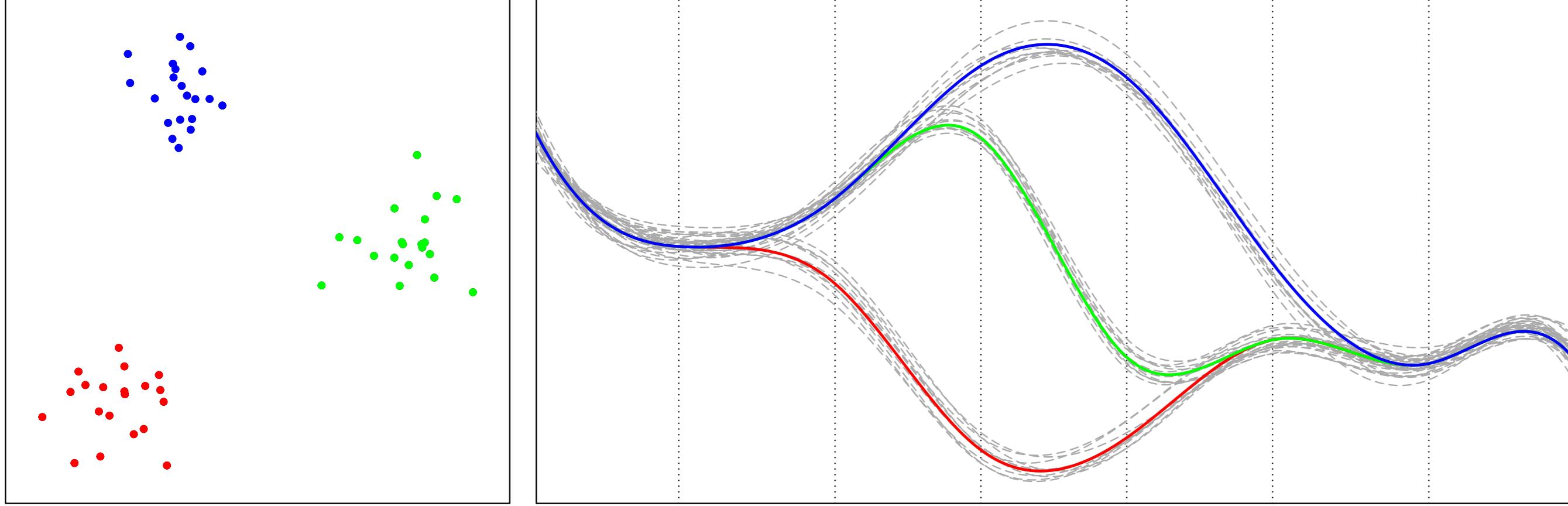


Figure 1: Comparison of a clustering problem for non-functional and functional data.

## Contributions

We propose a Bayesian approach for the analysis of functional data that

- exploits the local property of B-spline basis expansion to perform indirect local clustering,
- defines **unit-specific B-spline parameters in terms of unit- and basis-specific cluster assignments**,
- employs an **ad-hoc definition for contiguous cluster-specific parameters ensuring smooth functions**,
- employs a novel **dependent random partition model inducing sequences of random partitions exhibiting semi-Markovian dependence**.

## Modeling expected values via B-spline basis expansion

We assume that a random curve  $Y_i(x)$ , evaluated at the point  $x \in \mathbb{R}$ , follows

$$Y_i(x) \mid \boldsymbol{\theta}^*, \mathbf{c}_i, \sigma^2 \stackrel{ind}{\sim} \mathcal{N}(\mathbf{b}(x)^\top \boldsymbol{\theta}_i, \sigma^2),$$

where  $\mathbf{b}(x)$  is a  $d$ -degree B-spline with basis coefficient

$$\boldsymbol{\theta}_i = (\theta_{1,c_{i,1}}^*, \dots, \theta_{K,c_{i,K}}^*),$$

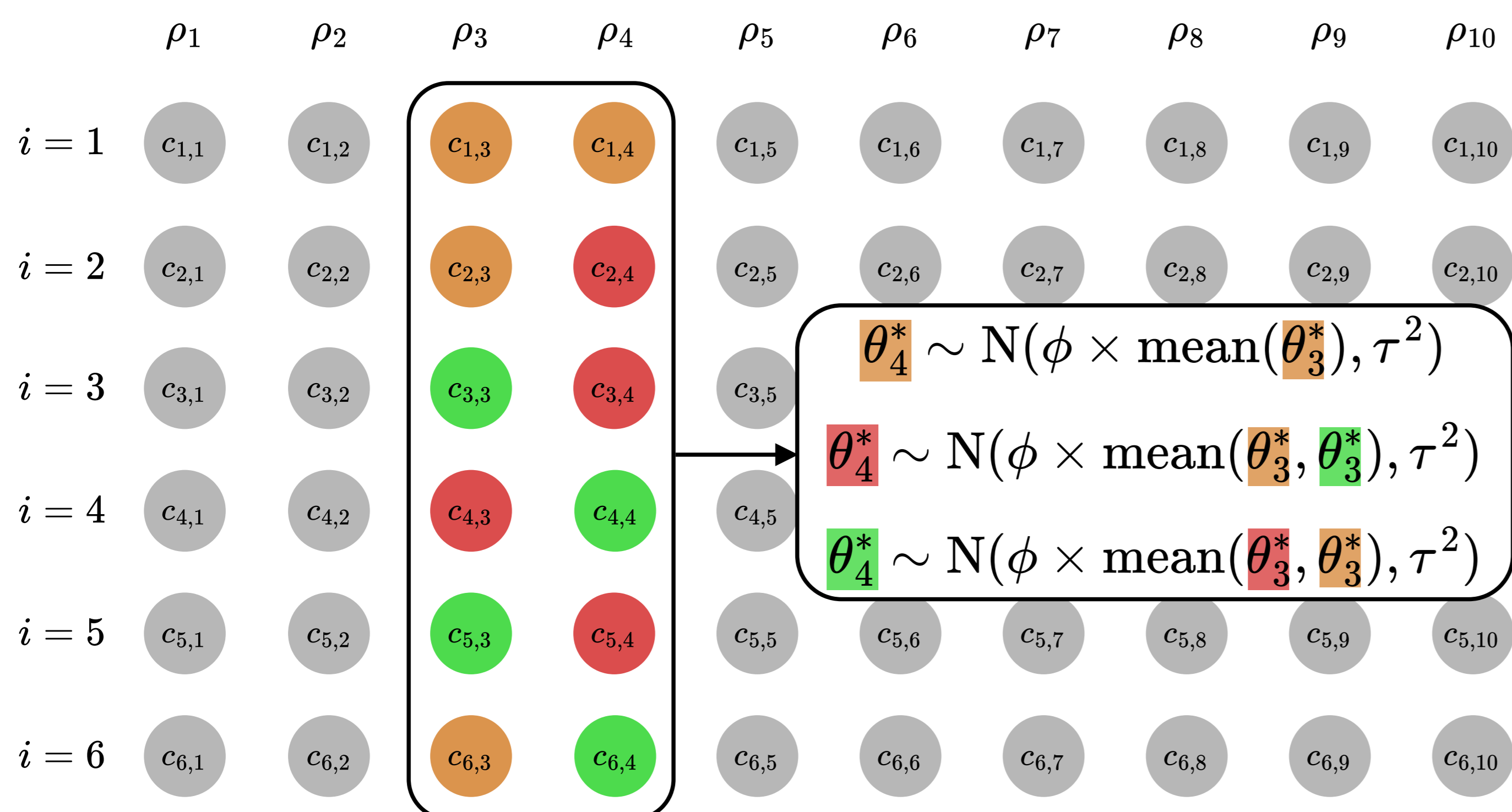
determined by  $K$  cluster assignments  $\mathbf{c}_i = (c_{i,1}, \dots, c_{i,K})$ , one for each basis.

**Local property:** the expected values of different functions,  $\mathbf{b}(x)^\top \boldsymbol{\theta}_i$ , coincide in part of their domain if the functions share local clusters for enough contiguous bases.

## Modeling cluster-specific parameters

We consider a partition of the curves at each basis,  $\rho_k$ , and define the cluster-specific parameters in such a way that contiguous parameters in each  $\boldsymbol{\theta}_i$  are similar:

$$\theta_{kj}^* \mid \boldsymbol{\theta}_{k-1}^*, \rho_k, \rho_{k-1}, \phi, \tau^2 \stackrel{ind}{\sim} \mathcal{N}\left(\frac{\phi}{|C_{k-1}^{(\rightarrow j)}|} \sum_{l \in C_{k-1}^{(\rightarrow j)}} \theta_{k-1,l}^*, \tau^2\right).$$



## References

- Page, G. L., Quintana, F. A., Dahl, D. B. (2022). Dependent Modeling of Temporal Sequences of Random Partitions. In: Journal of Computational and Graphical Statistics, **31** (2), 614–627.
- Hubert, L., and Arabie, P. (1985). Comparing Partitions. In: Journal of Classification, **2**, 193–218.
- Polson, N. G., Scott, J. G., Windle, J. (2013). Bayesian inference for logistic models using pólya–gamma latent variables. In: Journal of the American Statistical Association, **108** (504), 1339–1349..

## Semi-Markovian Random Partition Model (smRPM)

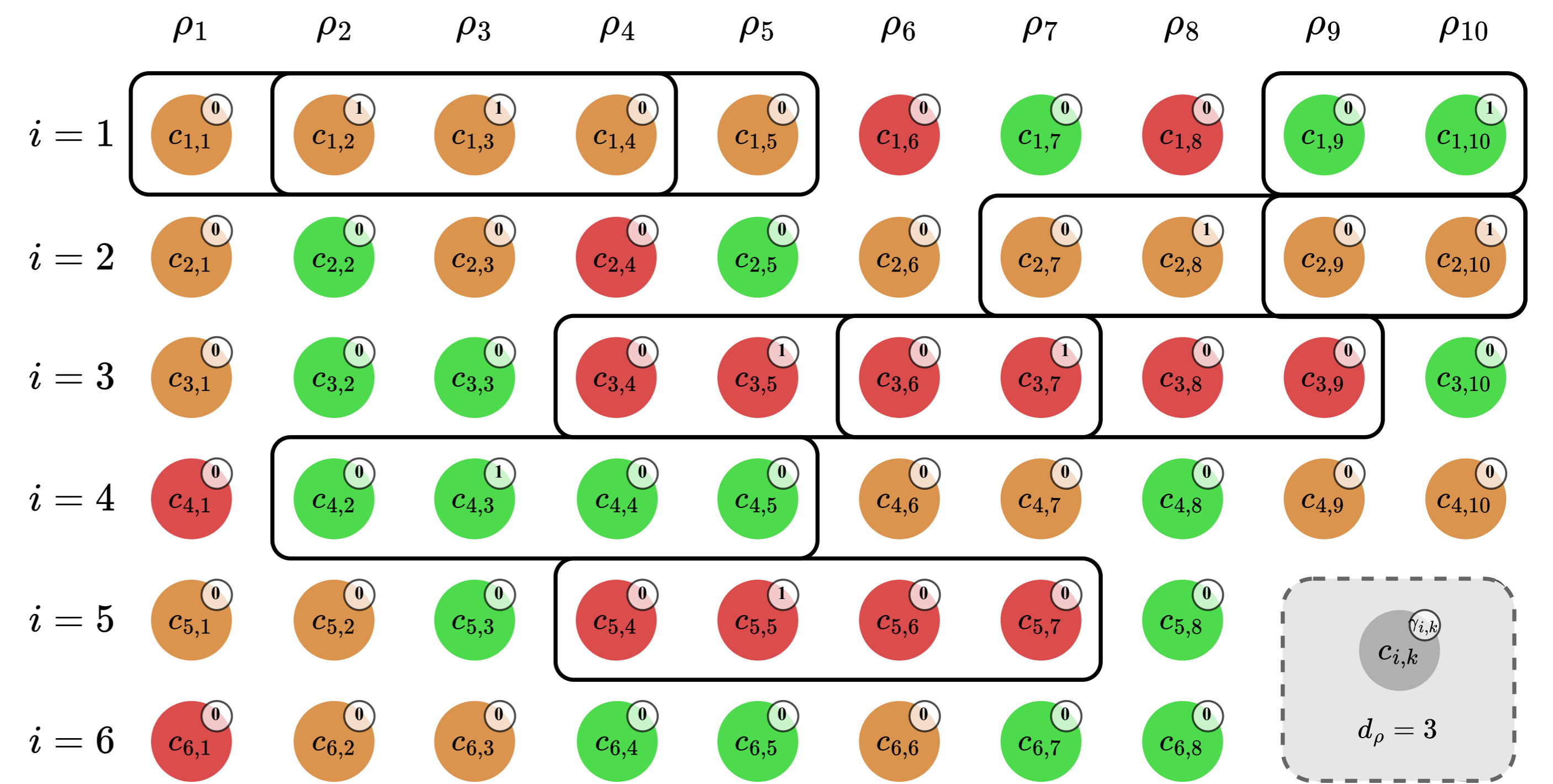
We introduce auxiliary variables explicitly modeling the evolution of the partitions at different basis:

$$\gamma_{ik} = \begin{cases} 1 & \text{if curve } i \text{ cannot be reallocated when moving from basis } k-1 \text{ to } k+d_\rho-1 \\ 0 & \text{otherwise} \end{cases}.$$

The distribution of the partition at basis  $k$  is influenced by the auxiliary variables related to bases  $k, \dots, k-d_\rho+1$ , collected in  $\boldsymbol{\gamma}_k^{(d_\rho)}$ :

$$\Pr(\rho_k = \lambda \mid \boldsymbol{\gamma}_k^{(d_\rho)}, \rho_{k-1}) = \frac{\Pr(\rho_k = \lambda) \mathbf{I}(\lambda \in P_{\mathcal{R}_k})}{\sum_{\lambda' \in P} \Pr(\rho_k = \lambda') \mathbf{I}(\lambda' \in P_{\mathcal{R}_k})},$$

where  $P_{\mathcal{R}_k}$  is the set of partitions that are compatible with  $\rho_{k-1}$  based on  $\boldsymbol{\gamma}_k^{(d_\rho)}$ .



The probabilities of success of the auxiliary variables at different basis can be independent,

$$\begin{aligned} \gamma_{ik} \mid \alpha_k &\stackrel{ind}{\sim} \text{Ber}(\alpha_k), \\ \alpha_k &\stackrel{iid}{\sim} \text{Beta}(a_\alpha, b_\alpha), \end{aligned}$$

or have  $d_\gamma$ -order dependence,

$$\begin{aligned} \gamma_{ik} \mid \boldsymbol{\alpha}, \boldsymbol{\gamma}_{i,k-1}^{(d_\gamma)} &\stackrel{ind}{\sim} \text{Ber}\left(\pi\left(\alpha_0 + \alpha_1 \sum_{q=1}^{d_\gamma} \gamma_{i,k-q}\right)\right), \\ \boldsymbol{\alpha} &\sim \mathcal{N}_2(\mathbf{a}, \mathbf{A}), \quad \omega_{ik} \stackrel{iid}{\sim} \text{PG}(0, 1). \end{aligned} \quad (\text{Polson et al. (2013)})$$

**Notation:** assuming  $\rho_1 \sim \text{CRP}(M)$ , the model is denoted as  $\text{smRPM}_{d_\rho, d_\gamma}(\boldsymbol{\alpha}, M)$ .

If  $d_\rho = 1, d_\gamma = 0$ , the model coincides with temporal Random Partition Model (Page et al., 2022).

## Simulations

We define  $n^{(ref)} = 5$  reference functional observations with known cluster assignments, and we simulate  $R = 50$  datasets containing  $n^{(rep)} \in \{10, 30\}$  realizations of each of these reference functional observations.

We simulate the observations and the cluster-specific parameters as

$$\begin{aligned} Y_i(x) \mid \boldsymbol{\theta}^*, \mathbf{c}_i, \sigma^2 &\stackrel{ind}{\sim} \mathcal{N}(\mathbf{b}(x)^\top \boldsymbol{\theta}_i, \sigma^2), \\ \theta_{1j}^* &\stackrel{ind}{\sim} (10j, 5), \quad \theta_{kj}^* \mid \boldsymbol{\theta}_{k-1}^*, \rho_k, \rho_{k-1} \stackrel{ind}{\sim} \left(\frac{1}{|C_{k-1}^{(\rightarrow j)}|} \sum_{l \in C_{k-1}^{(\rightarrow j)}} \theta_{k-1,l}^*, 5\right), \end{aligned}$$

where  $\sigma^2 \in \{1, 4\}$  quantifies the noise added to the reference functional observations.

**Evaluation:** Adjusted Rand Index (Hubert and Arabie, 1985) computed on each MCMC iteration.

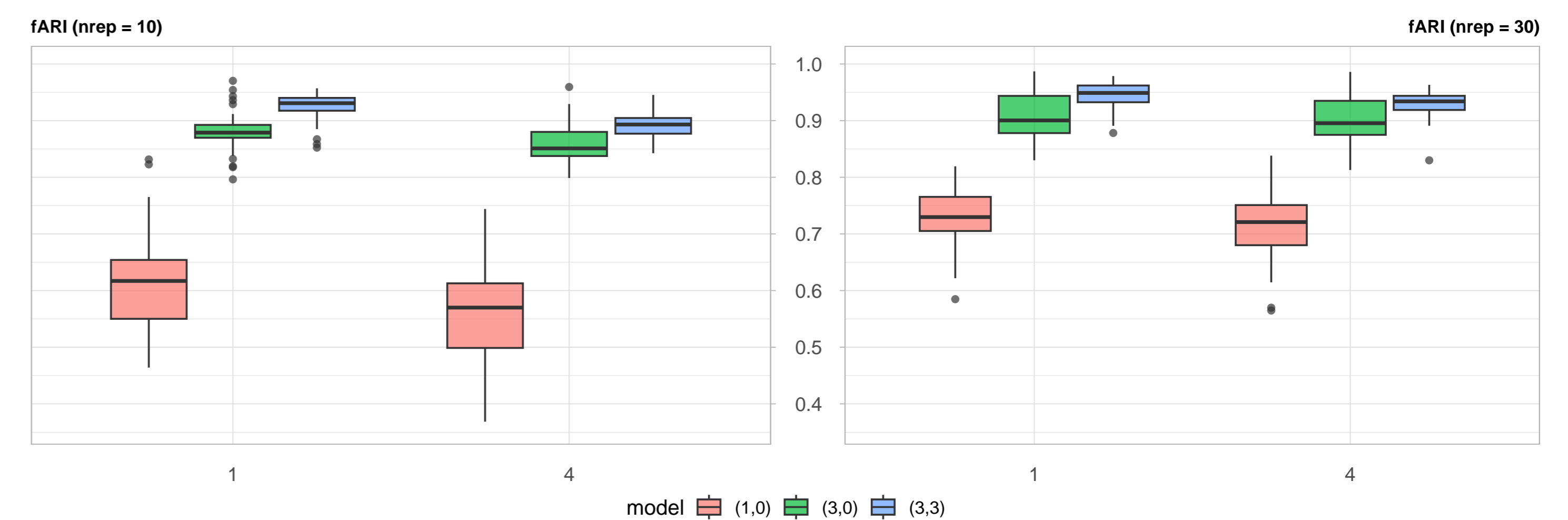


Figure 2: Boxplot of the average posterior ARI across the sequence of partitions computed on  $R = 50$  simulated datasets.