

# A Regression Mixture Model to understand the effect of the Covid-19 pandemic on Public Transport Ridership

Hugues Moreau  
Université Gustave Eiffel  
Paris, France

hugues.moreau.pro@gmail.com

Étienne Côme  
Université Gustave Eiffel  
Paris, France

etienne.come@univ-eiffel.fr

Allou Samé  
Université Gustave Eiffel  
Paris, France

allou.same@univ-eiffel.fr

Latifa Oukhellou  
Université Gustave Eiffel  
Paris, France

latifa.oukhellou@univ-eiffel.fr

**Abstract**—The Covid-19 pandemic drastically changed urban mobility, both during the height of the pandemic with government lockdowns, but also in the longer term with the adoption of working-from-home policies. To understand its effects on rail public transport ridership, we propose a dedicated Regression Mixture Model able to perform both the clustering of public transport stations and the segmentation of time periods, while ignoring variations due to additional variables such as the official lockdowns or non-working days. Each cluster is thus defined by a series of segments in which the effect of the exogenous variables is constant. As each segment within a cluster has its own regression coefficients to model the impact of the covariates, we analyze how these coefficients evolve to understand the changes in the cluster. We present the regression mixture model and the parameter estimation using the EM algorithm, before demonstrating the benefits of the model on both simulated and real data. Thanks to a five-year dataset of the ridership in the Paris public transport system, we analyze the impact of the pandemic, not only in terms of the number of travelers but also on the weekly commute. We further analyze the specific changes that the pandemic caused inside each cluster.

**Index Terms**—Clustering, Segmentation, Mixture of Regressions, Generative models, Expectation-Maximization

## I. INTRODUCTION

Finding similarities between groups of samples, a problem known as clustering, has been an active research domain in the last decades. However, when the samples are time series, one might want to find similarities between different periods. Segmentation, or change-point detection [19], is the problem of dividing the temporal axis into intervals in which the data follow a constant distribution. To combine these two approaches, we can use a *sequential* approach, by computing clusters first and finding change points inside each cluster. However, sequential approaches are noticeably unoptimized. Consequently, researchers have developed joint approaches, such as SegClust by Picard *et al.* [16], or ClustSeg [18].

The introduction of exogenous variables in the mixture model serves two purposes: firstly, it can help the model to remain unaffected by their effects. In the urban mobility domain, it could be interesting to include exogenous variables that represent the three lockdowns in order to enable the model to explain the changes in the data using regression

coefficients rather than new segments. Secondly we want our model to be able to use the influence of these variables to be important to our model: it is worth considering that during typical pandemic behavior, the disparity between weekends and weekdays may be smaller than in the period prior to the year 2020. This is the rationale of the publications which perform a joint regression and clustering [7] or regression and segmentation [11], [12]. By explaining the variations in the data with exogenous variables, adding regression to a clustering-segmentation model helps the segmentation to detect only change points for meaningful variations, instead of reacting to variations that can be considered trivial or easily explainable given the context.

For the analysis of public transport ridership, clustering and regression have already been used successfully either to know which types of travel need improvement [15], or to determine how investigate one metro line can help alleviate the flow of passengers when another line or station is unavailable [6]. However, these studies focus on signals for which the distribution has the same parameters for the whole duration. The goal in the present paper is to go further and include the effects of the pandemic in the analysis. For instance, to model the difference between weekdays and weekends, we may compute one mean for each type of day [15], or with an additive model [6]. The latter allows the model to combine the influence of several variables, such as the day of the week and the position in the year. Furthermore, having several segments, each with its regression coefficients, could help us understand how the pandemic changed travel behaviors by looking at which coefficients differ before and after the pandemic crisis.

This paper aims at better understanding how the pandemic affected the usage of the Paris public transport network. To do so, we used a public transport ridership dataset published by the transport organisation authority called Île-de-France Mobilités (IdFM) which contains records of public transport entries on the rail network of Paris city and its suburbs between 2017 and June 2022. The data cover the start of the pandemic, as well as the three lockdowns. Obviously, the official lockdowns had a strong impact on public travel but we are also interested in the impact of the pandemic on travel

behaviors. We would expect that, if some stations are relatively unaffected by the pandemic, a clustering-segmentation model could create a cluster for these stations, in which the transition points are located at different dates than the pandemic.

The present work focuses on the model that allows this interpretation to occur, and its contributions can be summarized as follows:

- we propose a regression mixture model which performs clustering and segmentation, and allows the data to follow a regression model where covariates can be of any type.
- we develop a formulation that models both univariate and multivariate data.
- we evaluate this model against a complete set of baselines, both on synthetic and real data.

The rest of this paper is organized as follows: after a literature review, section III presents both our model and the parameter estimation algorithm. Then, we evaluate our proposed model, on synthetic data (section IV-B) and on the public transport ridership dataset in section IV-D. Finally, section V provides some concluding remarks and proposes options for future work. To foster research on this domain, we publish our code at [https://github.com/HuguesMoreau/GMM\\_Clustering\\_Segmentation](https://github.com/HuguesMoreau/GMM_Clustering_Segmentation).

## II. RELATED WORK

Our work consists of the intersection of three research domains: clustering, segmentation, and regression.

As the scope of these domains is very broad, we provide only a concise overview, focusing on their intersections.

Clustering is a field of research aiming to group samples according to their similarity. We will use the formalism of *Gaussian Mixture Models*, it is assumed that the cluster each sample belongs to is a hidden variable, which is estimated using the *Expectation-Maximization* (EM) algorithm [8].

Like clustering, segmentation is also an active domain of research, that consists in finding intervals of time in which the behavior of the variables is as constant as possible. The temporal ordering allows the formulation of dynamic programming algorithms that guarantee to find a global optimum of the objective function (see Truong *et al.* [19] for a comprehensive review). Apart from dynamic programming algorithms and their variants, another type of model used to perform the segmentation of time series are Hidden Markov Models [10]. However, the transitions of Hidden Markov Models cannot always be controlled as precisely as required without additional constraints [14].

When the segmentation is paired with clustering, using dynamic programming requires re-computing a clustering for each time step (as in [16]) which might make the whole algorithm too long to compute. To prevent this problem, Samé *et al.* [18] developed a model named *ClustSeg* that performs both clustering and segmentation. To obtain an estimate in a reasonable time, the authors also present the rules to estimate the whole partition using a single EM run. This is the approach we expand by adding the contribution of exogenous variables. Note that the *ClustSeg* model already projects the data on a

basis of functions that can be understood as exogenous variables that are independent of the individuals. This is the reason why Chamroukhi renamed this model Piecewise Regression [3], even though the regression is, as in [18], only on temporal variables. Our contribution is to allow the incorporation of any type of variable and not only time covariates.

To account for the exogenous variables, we will use a linear model. De Veaux [7] combined linear regression with clustering to obtain a *mixture of regression*, a model where the mean of each cluster is a linear combination of some other variables. The author also proposed an adaptation of the EM algorithm to estimate the parameters of the model, which is still relatively close to the estimation procedure of classical mixture models. Another family of models that falls in the intersection between clustering and regression consists in using use a model named Latent class clustering-based random parameter ordered logit model (LCROL, [4]), where the probability for any individual to belong to a cluster follows a logistic regression, depending on the provided covariates. This idea is similar to the segmentation used by Samé *et al.* [18] and Chamrouki [3], except that the only variable available in the logistic regression is time (for further details, see section III-B).

The last field of interest is the intersection of segmentation and regression, named *segmented regression*. However, this name often designates studies where the breakpoint is fixed before the start of the analysis. This kind of segmented regression has been used in medicine [2], or for the evaluation of public policies [13]. Most of the time, only two segments are used to know whether the difference between the two is significant.

In our case, we will exhibit the results our model obtains on a public transport dataset. But first, we need to present the formulation of our model, along with the formulation of the EM algorithm to estimate the parameters of the mixture model.

## III. MODEL DEFINITION AND PARAMETER ESTIMATION

### A. Definitions and notations

The data to be analyzed in this article consist of a set of time series observed over the same time grid indexed by  $t$ .

The data we want to model are a set of *individuals*, which we want to group into  $K$  *clusters*. The dataset covers a time series observed over discrete timesteps  $t$  (in our case, days). For each individual  $i$ , and each day  $t$ , the data is a  $D$ -dimensional vector noted  $\mathbf{y}_{i,t}$ . Note that the model does not assume that all couples of individuals and days are present: some may be missing from the dataset (and will be in the case of the Public Transport dataset, see fig. 1). In all cases, we want to obtain, for each cluster  $k$ , a series of  $S$  *segments*, that is, an interval of timestamps with similar behavior. However, there are several causes for change: in the data which we want to exclude from our analysis. This is why we include *exogenous variables* that will help to explain the changes in the observed variables without resorting to new segments. We note  $x_{i,t,l}$  the value of the  $l^{\text{th}}$  variable for individual  $i$  and

timestep  $t$ . We assume that all variables are continuous: if one variable has more than two categories, we break it down into several dummy (zero-one) variables. Note that in practice, the variables we will use are mostly constant along all timesteps, or for all individuals.

### B. Model definition

Let  $z_i$  be the cluster assignment of individual  $i$ , and  $w_{i,t}$  the segment of individual  $i$  at time  $t$ . We assume that the cluster assignments  $z_i$  are i.i.d., following a multinomial distribution whose weights are parameters of our model:

$$P(z_i = k) = \pi_k \quad \text{s.t.} \quad \sum_{k=1}^K \pi_k = 1. \quad (1)$$

Given the cluster  $k$  of individual  $i$ , we assume that the segment assignment  $w_{i,t}$  follows a multinomial law defined by the probabilities:

$$P(w_{i,t} = s | z_i = k) = \kappa_t^{k,s} = \frac{\exp(t \cdot u_{k,s} + v_{k,s})}{\sum_{s'} \exp(t \cdot u_{k,s'} + v_{k,s'})}, \quad (2)$$

where  $u_{k,s}$  and  $v_{k,s}$  are model parameters. In other words, the only reason why we see definitive transitions is the fact that the arguments of the softmax functions are monotonous. This way of modelling segments using independent day assignments might seem unusual. We use the formulation in equation 2 because the independence of the segment assignments allows us to make use of the EM algorithm (see section III-C). The partition of the time axis resulting from the application of the *Maximum a Posteriori* (MAP) rule may not be a series of contiguous segments. We thus applied the following constraints to make sure that the distributions of the temporal class  $\kappa_t^{k,s}$  vary faster than a given threshold  $\lambda$ :

$$\forall k, \forall s \in \{1 \dots S - 1\}, \quad u_{k,s+1} - u_{k,s} > \lambda \quad (3)$$

Note that  $\lambda$  is not a parameter of our model, but a value we choose before any parameter estimation. Numerically, we set it in such a way that the temporal classes distributions  $\kappa_t^{k,s}$  go from 0.01 to 0.99 in less than three months<sup>1</sup>. We found these constraints paramount: without them, the distribution of temporal classes  $\kappa_t^{k,s}$  would remain nearly constant (equal to approximately  $1/S$ , where  $S$  is the number of segments), and the posteriors would not have the contiguity expected from segments.

*Observation model:* Similarly to ClustSeg [18], our observation model is a Gaussian linear regression model:

$$\mathbf{y}_{i,t} | (z_i = k, w_{i,t} = s) \sim \mathcal{N}(\boldsymbol{\mu}_{i,t}^{k,s}, \boldsymbol{\Sigma}_{k,s}) \quad (4)$$

However, contrary to ClustSeg, in which the statistical mean  $\boldsymbol{\mu}$  depended only on the cluster, segment, and time covariates,

<sup>1</sup>To do so, we solve  $\exp(\lambda \cdot t + c) = 0.01$  and  $\exp(\lambda(t + \Delta t) + c) = 0.99$ , yielding  $\lambda = (\log(0.99/0.01))/\Delta t$ . As the time  $t$  varies between zero and one in all our formulas, the corresponding  $\Delta t$  for three months is equal to 90 days divided by the number of days in the interval  $T$ .

our formulation introduces exogenous variables that vary for both individuals  $i$  and days  $t$ . As in mixtures of regressions [7], the mean is a sum of the linear contributions of all exogenous variables:

$$\boldsymbol{\mu}_{i,t}^{k,s} = \mathbf{m}_{k,s} + \sum_{l=1}^L x_{i,t,l} \boldsymbol{\alpha}_{l,k,s}, \quad (5)$$

where  $\mathbf{m}_{k,s}$  is the intercept of the linear regression model, and the  $\boldsymbol{\alpha}_{l,k,s}$  are the contributions of the  $L$  temporal variables.

For the covariance matrix of cluster  $k$ , segment  $s$ ,  $\boldsymbol{\Sigma}_{k,s}$ , we choose to have a diagonal covariance (meaning that all dimensions are independent) but other constraints are possible (full covariance, spherical, *etc.*), similarly to the fact that one can choose the type of covariance for classical Gaussian Mixture Models.

To summarize this section, the free parameters of the model, denoted as  $\boldsymbol{\theta}$ , are the mixture proportions  $(\pi_k)_k$ , the segmentation parameters  $(u_{k,s}, v_{k,s})_{k,s}$ , the mean  $\mathbf{m}_{k,s}$  of cluster  $k$  and segment  $s$ , the contributions of the exogenous variables  $(\boldsymbol{\alpha}_{l,k,s})_{l,k,s}$ , and the covariance matrices  $(\boldsymbol{\Sigma}_{k,s})_{k,s}$ . The next section explains how to estimate them via the EM algorithm [8].

### C. Parameter estimation

The parameter estimation is done by maximizing the log-likelihood:

$$\begin{aligned} \mathcal{L}(\boldsymbol{\theta}) &= \log P(\mathbf{Y} | \mathbf{X}, \boldsymbol{\theta}) \\ &= \sum_i \log \sum_k \prod_t \pi_k \sum_s \kappa_t^{k,s} \mathcal{N}(\mathbf{y}_{i,t}; \boldsymbol{\mu}_{i,t}^{k,s}, \boldsymbol{\Sigma}_{k,s}) \end{aligned} \quad (6)$$

The formulation of our model allows us to use the Expectation-Maximization algorithm [8], a powerful algorithm aiming to estimate the parameters of latent (hidden) variable models. In our case, the latent variable is the union of the cluster and segment, which we estimate using both the observed variable  $(\mathbf{y}_{i,t})$  and the exogenous variables  $(\mathbf{x}_{i,t})$ . The EM algorithm consists in alternating two steps.

*E-step:* In this step, we compute the posterior distribution of the latent variables, using the values of the parameters at the current iteration. As mentioned in [18], this is done in two steps. First, we need to compute the posterior probability of a given individual  $i$  to belong to cluster  $k$ , which we note  $\rho_i^k$ :

$$\rho_i^k = P(z_i = k | \mathbf{y}_i, \mathbf{x}_i; \boldsymbol{\theta}) \quad (7)$$

$$= \frac{\pi_k \prod_t \sum_s \kappa_t^{k,s} \mathcal{N}(\mathbf{y}_{i,t}; \boldsymbol{\mu}_{i,t}^{k,s}, \boldsymbol{\Sigma}_{k,s})}{\sum_{k'} \pi_{k'} \prod_t \sum_s \kappa_t^{k',s} \mathcal{N}(\mathbf{y}_{i,t}; \boldsymbol{\mu}_{i,t}^{k',s}, \boldsymbol{\Sigma}_{k',s})}, \quad (8)$$

where we note  $\mathbf{y}_i = (\mathbf{y}_{i,t})_{t \in \{1, \dots, T\}}$ .

Then, we can obtain the posterior probability for any day  $t$  to belong to segment  $s$ :

$$r_{i,t}^{k,s} = P(w_{i,t} = s, z_i = k | \mathbf{y}_i, \mathbf{x}_i; \boldsymbol{\theta}) \quad (9)$$

$$= \rho_i^k \frac{\kappa_t^{k,s} \mathcal{N}(\mathbf{y}_{i,t}; \boldsymbol{\mu}_{i,t}^{k,s}, \boldsymbol{\Sigma}_{k,s})}{\sum_{s'} \kappa_t^{k,s'} \mathcal{N}(\mathbf{y}_{i,t}; \boldsymbol{\mu}_{i,t}^{k,s'}, \boldsymbol{\Sigma}_{k,s'})} \quad (10)$$

$$r_{i,t}^{k,s} = P(z_i = k, w_{i,t} = s | \mathbf{y}_i, \mathbf{x}_i; \boldsymbol{\theta}), \quad (11)$$

Once we have computed the responsibilities, we update the parameters during the *Maximization* step.

*M-step:* In the maximization step, we use the responsibilities  $r_{i,t}^{k,s}$  computed earlier, and obtain the parameters that maximize the expectation of the log-likelihood. Conveniently, we can write it as a sum of terms which can be optimized separately:

$$\begin{aligned} \mathcal{Q}(\mathbf{Y}, \mathbf{x}; \boldsymbol{\theta}) &= \sum_{i,t,k,s} r_{i,t}^{k,s} \log \pi_k \\ &+ \sum_{i,t,k,s} r_{i,t}^{k,s} \log \frac{\exp(t.u_{k,s} + v_{k,s})}{\sum_{s'} \exp(t.u_{k,s'} + v_{k,s'})} \quad (12) \\ &+ \sum_{i,t,k,s} r_{i,t}^{k,s} \log \mathcal{N}(\mathbf{y}_{i,t}; \boldsymbol{\mu}_{i,t}^{k,s}, \boldsymbol{\Sigma}_{k,s}) \end{aligned}$$

Thus, each set of parameter values can be computed independently from the others, the only exception being the variance which requires the segment mean and variable contributions of the current step. In this section, we display only the results of the maximization step.

*a) cluster proportions:* As with many mixture models, the mixture proportions are the proportion of the responsibilities of the samples. For each cluster  $k$ :

$$\pi_k = \left( \sum_{i,t,s} r_{i,t}^{k,s} \right) / \left( \sum_{i,t,k',s} r_{i,t}^{k',s} \right) \quad (13)$$

*b) segment borders:* The segment parameters  $u, v$  are found by computing the maximum of the following logistic regression problem:

$$\sum_{i,t,k,s} r_{i,t}^{k,s} \log \frac{\exp(t.u_{k,s} + v_{k,s})}{\sum_{s'} \exp(t.u_{k,s'} + v_{k,s'})} \quad (14)$$

This problem is convex (even with the constraints mentioned in equation 3), and we solve it using `scipy`'s `optimize` package implementation of the quasi-Newton method from [5]. We ensure identifiability by setting  $\sum_s u_{k,s} = \sum_s v_{k,s} = 0$  for all  $k$ .

*c) mean, regression coefficients, and covariance:* We can find the optimal values for the last term of equation 12 separately for each segment. For all  $k, s$ , the values of  $\mathbf{m}_{k,s}$ ,  $\boldsymbol{\alpha}_{k,s}$ , and  $\boldsymbol{\Sigma}_{k,s}$  minimize:

$$\sum_{i,t} r_{i,t}^{k,s} \left( (\mathbf{y}_{i,t} - \boldsymbol{\mu}_{i,t}^{k,s}) \boldsymbol{\Sigma}_{k,s}^{-1} (\mathbf{y}_{i,t} - \boldsymbol{\mu}_{i,t}^{k,s})' + \log \det \boldsymbol{\Sigma}_{k,s} \right), \quad (15)$$

where  $\mathbf{x}'$  denotes the transpose of vector  $\mathbf{x}$  and  $\boldsymbol{\mu}_{i,t}^{k,s} = \mathbf{m}_{k,s} + \sum_l x_{i,t,l} \cdot \boldsymbol{\alpha}_{l,k,s}$ . This is a common linear regression problem where the  $r_{i,t}^{k,s}$  are weights [17].

As mentioned earlier, we use these formulas during each of the maximization steps to estimate new values of the parameters. The complete algorithm can be summarized as follows:

---

**Algorithm 1** the EM algorithm

---

**Require:** observed data  $(\mathbf{y}_{i,t})_{i,t}$ , exogenous variables  $(\mathbf{x}_{i,t})_{i,t}$ , a number of clusters  $K$  and segments  $S$ , and an initial parameter  $\boldsymbol{\theta}$

**while** the log-likelihood has not converged **do**

*E-step:* Compute the expectation of the latent variables  $r_{i,t}^{k,s}$  using the parameters of the model  $\boldsymbol{\theta}$

*M-step:* Find the parameters of the model  $\boldsymbol{\theta}$  that best explain the data

**end while**

---

We keep on alternating the expectation step and maximization steps until the log-likelihood stops increasing by more than  $10^{-4}$  in ten iterations. Only then is parameter estimation over and we can start evaluating our model.

*Initialization:* To initialize our algorithm, we start by setting the mixture proportions equal to  $1/K$ , where  $K$  is the number of clusters. The segment parameters  $u, v$  are set in such a way as to respect the constraints we set (see equation 3), and the borders between segments delimit segments of equal duration. The segment means  $\mathbf{m}_{k,s}$  are estimated by assigning each individual a random cluster and each day its corresponding segment, and computing the mean of values inside each segment. Finally, we end the initialization by setting the contribution of exogenous variables to zero.

## IV. EXPERIMENTS

This section is devoted to presenting the experimental results that highlight the efficacy of the proposed model. We will compare the performance of our model against a set of baselines using artificial data. Subsequently, we will assess its performance on the public transport ridership dataset through cross-validation. Lastly, we will analyze the results obtained by the model to gain insights into the stations that experienced the most significant impacts due to behavioral changes resulting from the pandemic.

### A. Protocol

The model presented above allows us to cover clustering, segmentation, and regression at once. To justify its use despite the additional complexity, we compare it to several alternatives. Most of these methods consist in using ClustSeg [18] without exogenous variables, which means that the basis of constant functions is used to represent the time series. Conversely, as the proposed model includes covariates in the parameter estimation, the basis of functions we choose to represent the time series can be found in the variables we have (in our case, splines of degree 2). To demonstrate the

interest of simultaneously performing the segmentation and the clustering, we include in our comparison two methods that compute clusters before segments:

- $\text{Reg} \rightarrow \text{Clust} \rightarrow \text{Seg}$ : This method operates in three sequential steps. Firstly, we estimate regression coefficients ( $\alpha$ , see eq. 5) common to the whole dataset. Secondly, we perform a clustering of the residual time series derived from the first step. Thirdly, we perform a segmentation of the clusters obtained in the second step. Note that for synthetic data, we generate the ground truth regression coefficients with uniform distributions. Hence, computing regression coefficients on the whole dataset does not make much sense. In other words, we expect this method to be completely irrelevant for synthetic data. However, it might prove somewhat useful for real data, when the contributions of the exogenous variables might be similar across clusters and segments. We include it nonetheless for both types of data.
- $(\text{Clust}+\text{Reg}) \rightarrow (\text{Seg}+\text{Reg})$ : This baseline involves fitting a mixture of regression model, before performing a segmented regression inside each cluster. It is important to note that we estimate new variable contributions during both the clustering and segmentation processes.

We also include three baselines where the clusters are estimated at the same time as the segments, meaning that the model can use the change points in the samples to estimate how close all individuals are to each cluster:

- $\text{Clust}+\text{Seg}$ : We simply reuse the method from [18], without any covariates. As  $\text{ClustSeg}$  requires a basis of (temporal) functions, using it without any covariates is equivalent to determining piecewise constant cluster prototypes. The rationale for including this method is mainly to demonstrate the relevance of including the exogenous variables with real data.
- $\text{Reg} \rightarrow (\text{Clust}+\text{Seg})$ : This model is similar to the one we propose, except that we assume that the linear regression parameters are constant across all clusters and segments. Note that to estimate the parameters for this model, we perform a linear regression on the whole dataset (before any affectation onto clusters and segments), before using  $\text{ClustSeg}$  (without any exogenous variables) on the residuals. This estimation procedure is equivalent to estimating the parameters of the model using the classical EM algorithm based on Gaussian mixture models. Similarly to  $\text{Reg} \rightarrow \text{Clust} \rightarrow \text{Seg}$ , we expect the common regression coefficients to be ill-adapted to synthetic data.
- $(\text{Clust}+\text{Seg}) \rightarrow \text{Reg}$ : This method involves first using  $\text{ClustSeg}$  [18] without any covariables, then performing a regression on each segment. It has the same set of parameters as the one we propose, the only difference between the two being that the model we present has access to changing regression coefficients during the clustering and segmentation. Given that both this method

and the first one ( $\text{Clust}+\text{Seg}$ ) compute the clusters and segments without any variables, they assign the same partition to the data, which is why we do not use it on synthetic data, where the evaluation criterion is the relevance of the partition.

Using different evaluation criteria, we compare the partitioning obtained with our model to the ones resulting from these methods. For synthetic data (simulations), the creation of the dataset allows us to have access to the latent variables that generated the data as ground truth (cluster and segment of each individual and day). We consider that two samples belong to the same partition if and only if their cluster *and* segment both match, and use the Adjusted Rand Score to compare to the evaluated models' estimations.

For real-life data, there is no such reference to use. Hence, we look at the log-likelihood on an unseen dataset: we separate the dataset into two subsets for training and validation. We use the training dataset to estimate the parameters of the model. Then, we calculate the log-likelihood of the model by applying it to the validation dataset. The higher the log-likelihood, the better the model explains the validation data, and the better we consider our model.

### B. Validation of the model on synthetic data

*Parameter generation:* Before generating data, we need to fix the model parameters. We select arbitrarily  $K = 4$  clusters and  $S = 4$  segments. The mixture proportions  $\pi_k$  are drawn from a Dirichlet distribution with parameter 2, to make sure all clusters appear for at least one sample. Similarly, for each cluster, we draw segment proportions with this same distribution, where all slopes  $u$  are set so that a probability changes from 10% to 90% in one-tenth of the time interval. We sample the cluster and segment means  $m_{k,s}$  with  $\mathcal{N}(0, 1)$ . To imitate the exogenous variables from the real data, we use three exogenous variables: the first is constant across timesteps ( $x_{i,t,1} = y_{i,1,1}$ ), the second is constant across individuals ( $x_{i,t,2} = x_{1,t,2}$ ), and the third varies for both individuals and timesteps. All are drawn from  $\mathcal{N}(0, 1)$ . Then, we sample the contributions corresponding to these exogenous variables ( $\alpha$  coefficients) from  $\mathcal{N}(0, \Sigma_\alpha)$ , where the value of  $\Sigma_\alpha$  changes between experiments. Finally, we generate the observed variables using the variance  $\Sigma_{k,s} = 1$ . The choice of a relatively high unexplained variance  $\Sigma_{k,s}$  is voluntary, to place the model in a difficult scenario.

During the first set of experiments, we set  $\Sigma_\alpha = I$ , and change the number of individuals and timesteps in  $\{50, 100, 500, 1000\}$ . In the second series of tests, the number of individuals and timesteps remains constant and equal to 100 while the parameter that gives the importance of the contributions  $\Sigma_\alpha$  varies in  $\{0, 0.5, 1.0, 1.5\}$ . The case  $\Sigma_\alpha = 0$ , in particular, is interesting, for it means that the model will see variables that do not contribute anything to the observed variables. In this case, if its performance is significantly lower than a model that does not have access to the variables, this would mean that the proposed model overfits the exogenous variables. In order to mitigate the randomness, we repeat the

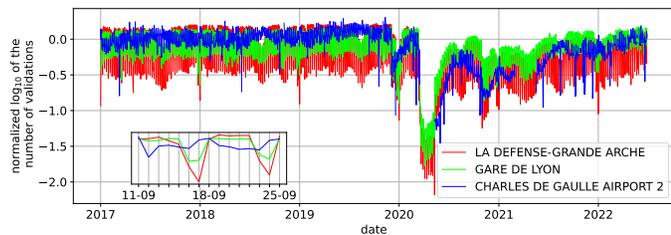


Fig. 1: Some examples of series in the dataset we process. The zoom-in represents two weeks starting September 11, 2019 (on a Monday).

experiments ten times for each combination of parameters, generating a whole new set of parameters each time.

*Results:* The results are summarized in table I. As expected, the models which do not take into account the variables (Clust+Seg) or which try to fit a regression before any clustering (Reg  $\rightarrow$  Clust  $\rightarrow$  Seg and Reg  $\rightarrow$  (Clust+Seg)) obtain ARI scores that are noticeably low. This is because the contributions we generated are independent of each other, meaning that fitting regression coefficients on the whole dataset always produces statistically insignificant parameters.

The only case when these models are relevant is the case where  $\Sigma_\alpha = 0$ , meaning that the exogenous contributions are always zero. For this experiment, the score only depends on whether clustering and segmentation are simultaneously carried out: Reg  $\rightarrow$  Clust  $\rightarrow$  Seg and (Clust+Reg)  $\rightarrow$  (Seg+Reg) are both worse than the other three models, which perform similarly. As expected, the position of the linear regression in the estimation process has no significant importance for this case.

### C. The Public Transport Ridership Dataset

To validate our model on real data, we used a recording of daily entries in the Parisian public transport system.<sup>2</sup> Gathered by the transport organization authority called Île-De-France Mobilités (IdFM), this dataset covers five and a half years, from January 1, 2017, to June 30, 2022. We also focus on the rail network (underground, public train, and tramway). Note that the records only cover the entries of travelers entering the network, which means that we have the origin of trips but not the destination. Note also that due to the organization of the ticketing system, one person may validate several times for a single trip when changing fare zones (if the journey includes connections between different modes). We do not try to compensate for this. Finally, we must emphasize that even though we use the term "individual" to designate each time series, one series of entries corresponds to a station, and not to a traveler.

<sup>2</sup>available at <https://data.iledefrance-mobilites.fr/explore/dataset/histo-validations-reseau-ferre/> and <https://data.iledefrance-mobilites.fr/explore/dataset/validations-reseau-ferre-nombre-validations-par-jour-1er-semester/information/>

a) *Preprocessing:* The raw data contain the number of entries for each type of ticket or subscription the travelers used. As we do not make use of this information, we simply consider the sum of all entries for each day and each station. However, only the couples (*station, day*) with at least one entry are present. One might be tempted to say that an absence of data for a given day means that no user entered this station (which might have happened during the lockdowns, for instance), but we have no way of knowing which days are missing due to actual errors in the data collection process, and which days are absent because of the lack of users. Thus, we leave missing data as-is.

To ensure anonymity, when a station recorded less than five entries for the same type of ticketing in a day, the dataset only contains the mention "less than five" for this type of ticket. In such cases, we considered that there were three entries for this day, station, and ticket type. We removed stations with more than 60% missing days or less than 500 entries per day on average. To handle the many noisy days in the data, we compute a two-week moving average of each station's series of entries. Each value that is below one-tenth of the result is removed.

Finally, given the strong imbalance between stations, we normalized the number of entries of each station. For each station, we divide the number of entries by the average number of entries during the first year, strikes excluded (from January 1, 2017 to December 31, 2017). We consider this period to be representative of usual behavior. Then, similarly to [6], we took the log in base 10 of the number of entries.

b) *Exogenous variables:* We include several dummy (binary) variables that relate to the type of day: one variable to know whether the day is a working day (meaning no holidays or weekends), seven variables for the days of the week, and two variables for strikes: as one strike was particularly between December 2019 and January 2020, we dedicate one variable to this period. The other 'strike' variable is the same for all other strikes. Similarly, we create three variables for the three lockdowns declared by the government. Finally, to account for tendencies throughout the year, such as a lower ridership during the summer, we use 20 year-periodic splines with degree 2. We also include four variables that are proper to each station, denoting the presence of four commercial types of train: underground (Paris only), RER (express regional network), transilien (a network going further than RERs), and commercial train (presence of national line in the same station).

### D. Experiments with real data

Section IV-B demonstrated that the model is able to find the true cluster of individuals generated using Gaussian distributions. This section shows that even with real data, the model we propose outperforms the baselines defined in section IV-B.

However, for many real-life datasets, we have no way of knowing which individuals belong to each cluster, or even if there is a 'right' number of clusters. We first begin to select the number of clusters using the slope heuristic [1].

	$\Sigma_\alpha = I$			
	$I = T = 50$	$I = T = 100$	$I = T = 500$	$I = T = 1000$
Reg $\rightarrow$ Clust $\rightarrow$ Seg	0.334 $\pm$ 0.152	0.293 $\pm$ 0.062	0.286 $\pm$ 0.061	0.377 $\pm$ 0.088
(Clust+Reg) $\rightarrow$ (Seg+Reg)	0.589 $\pm$ 0.141	0.643 $\pm$ 0.161	0.620 $\pm$ 0.081	0.788 $\pm$ 0.114
Reg $\rightarrow$ (Clust+Seg)	0.463 $\pm$ 0.099	0.440 $\pm$ 0.125	0.414 $\pm$ 0.112	0.511 $\pm$ 0.123
Clust+Seg	0.478 $\pm$ 0.106	0.453 $\pm$ 0.120	0.431 $\pm$ 0.120	0.561 $\pm$ 0.152
Clust+Seg+Reg (Proposed)	<b>0.708 <math>\pm</math> 0.143</b>	<b>0.784 <math>\pm</math> 0.110</b>	<b>0.865 <math>\pm</math> 0.088</b>	<b>0.895 <math>\pm</math> 0.072</b>

	$I = T = 100$			
	$\Sigma_\alpha = 0$	$\Sigma_\alpha = (0.5^2)I$	$\Sigma_\alpha = I$	$\Sigma_\alpha = (1.5^2)I$
Reg $\rightarrow$ Clust $\rightarrow$ Seg	0.568 $\pm$ 0.178	0.456 $\pm$ 0.125	0.346 $\pm$ 0.063	0.323 $\pm$ 0.090
(Clust+Reg) $\rightarrow$ (Seg+Reg)	0.542 $\pm$ 0.177	0.605 $\pm$ 0.122	0.661 $\pm$ 0.174	0.626 $\pm$ 0.145
Reg $\rightarrow$ (Clust+Seg)	0.645 $\pm$ 0.165	0.548 $\pm$ 0.094	0.467 $\pm$ 0.121	0.442 $\pm$ 0.133
Clust+Seg	0.686 $\pm$ 0.092	0.586 $\pm$ 0.137	0.499 $\pm$ 0.179	0.434 $\pm$ 0.107
Clust+Seg+Reg (Proposed)	<b>0.689 <math>\pm</math> 0.104</b>	<b>0.794 <math>\pm</math> 0.075</b>	<b>0.815 <math>\pm</math> 0.088</b>	<b>0.804 <math>\pm</math> 0.121</b>

TABLE I: The mean and standard deviation of the Adjusted Rand Score obtained with ten generations of synthetic data. The method (Clust+Seg)  $\rightarrow$  Reg is absent because it produces the same partition as Clust+Seg. The best result each time is in bold.

number of individuals	542 stations
days	2,004
Measures per day	1
covariates (number of variables)	type of transport available (4), day of the week (7), lockdowns (3), strikes (2), holidays (2), year-periodic splines (20), time (1)

TABLE II: An overview of the real dataset.

As we expect the data to exhibit one change point at the beginning of the pandemic, we select  $S = 2$  segments *a priori*. We begin by fitting the parameters for models with one to nineteen clusters. Each time, we obtain a (training) log-likelihood, which increases with the number of parameters. The slope heuristic consists in applying a penalty equal to twice the slope of the curve giving the log-likelihood as a function of the number of parameters in each model [1]. In our case, with  $S = 2$  segments, the optimal value is reached for  $K = 5$  clusters (results not shown).

To evaluate our model on real data, we perform five-fold cross-validation, dividing the dataset into five subsets of mutually-independent samples. This constraint of mutual independence is not trivial: if one split the dataset into separate periods, for instance, the folds would not be independent. The reason is that the cluster assigned to each individual remains constant: if we know the segment a given day belongs to, we can infer the individual’s cluster, which provides information about the cluster the individual belongs to outside of the training interval. This is why we perform cross-validation by dividing the dataset into five groups of individuals. We use the first four to estimate the parameters of the model, keeping the last group to measure the log-likelihood of the unseen data. We repeat the process four additional times, changing the validation dataset, and display the mean and standard deviation in table III.

We can draw several conclusions from the results in table III: firstly, the two methods that perform clustering and segmentation separately are noticeably worse than the rest which means that performing clustering and segmentation at the same time is paramount. Secondly, the model that imposes the regression coefficients to be equal across clusters and segments (Reg  $\rightarrow$  (Clust+Seg)) explains unseen data as well as a model that does not have access to the exogenous variables (Clust+Seg), hinting at overfitting.

Finally, the best two models are the ones that both perform clustering and segmentation at the same time, while still allowing variable contributions to differ between segments ((Clust+Seg)  $\rightarrow$  Reg and Clust+Seg+Reg). Between them, the proposed approach that performs clustering/segmentation and regression at the same time (Clust+Seg+Reg) is significantly better than the sequential model where the estimation of the contributions is done after the estimation of the partition ( $p < 0.01$  using an unpaired t-test).

However, this improvement comes at the cost of a higher computational burden: using a desktop computer<sup>3</sup>, the implementation that estimates all parameters at the same time with the EM algorithm takes fifteen minutes to converge on average ( $15.4 \pm 2.25$  min), which is four times longer to run than its sequential counterpart ( $4.0 \pm 0.6$  min for (Clust+Seg)  $\rightarrow$  Reg).

Using cross-validation, we demonstrated that the proposed model is able to reach a higher log-likelihood than its counterparts, indicating a better ability to model the data. We will now use this model to understand the effects of the pandemic on public transport ridership.

<sup>3</sup>a computer with an Intel I7 @ 2.80GHz and a 15 Go RAM, along with Ubuntu 20.04

Model	validation LL	number of parameters	Segments and clusters are computed ...	Regression coefficients are...
Reg $\rightarrow$ Clust $\rightarrow$ Seg	$-2.02 \times 10^6 \pm 3.0 \times 10^3$	73	Sequentially	Common
(Clust+Reg) $\rightarrow$ (Seg+Reg)	$-1.96 \times 10^6 \pm 8.2 \times 10^2$	424	Sequentially	Different
Clust+Seg	$2.03 \times 10^6 \pm 5.6 \times 10^3$	34	Simultaneously	Absent
Reg $\rightarrow$ (Clust+Seg)	$2.05 \times 10^6 \pm 5.9 \times 10^3$	73	Simultaneously	Common
(Clust+Seg) $\rightarrow$ Reg	$2.86 \times 10^6 \pm 7.0 \times 10^3$	424	Simultaneously	Different
Clust+Seg+Reg (Proposed)	<b><math>3.14 \times 10^6 \pm 8.2 \times 10^3</math></b>	424	Simultaneously	Different

TABLE III: The mean and standard deviation of the log-likelihood computed on the validation dataset, using five-fold cross-validation. The best result is in bold.

### E. Analysis of the results on public transport ridership

As mentioned at the beginning of the previous section, we apply the proposed model with five clusters and two segments. As the changes between segments are easier to explain than the difference between clusters, we will begin by explaining how the model provides insights into the public transport ridership before and after the pandemic, before using this behavior to understand how clusters differ from each other.

Figure 2 displays the sum of the mean and contributions of the time and the year-periodic splines. As expected, the mean of the second segment is noticeably lower than the first segment's, and the (positive) regression coefficient associated with the time models the recovery of each cluster. The impact of the pandemic on ridership is the gap between the reconstructions of the two segments, when their prior probability is equal (approximately early 2020, depending on the clusters). As expected, the gap between the extension of the first segment and the second segment is at its highest at the start of the pandemic, after which the number of entries increase to model the recovery of the network.

Given the importance of the pandemic on the population's trips, one could expect the model to pick up the start of the pandemic (and the first lockdown) as the change point between segments. What happened is slightly different: a major strike occurred between December 2019 and January 2020. Even though we created a dummy variable to model the effect of this specific strike, one variable is not enough to model the continuous recovery of the signal. As a consequence, for some clusters, the model underestimates the effect of this strike, and assigns the lowest days to the next segment, the segment dedicated to the pandemic. This is why the transition between segments is sometimes located at the very beginning of the year 2020.

The small number of days assigned this way to the second segment is not, however, enough to prevent the model from accounting correctly for the events of the pandemic. In March 2020, the French government decreed a lockdown to fight against the pandemic. While people were allowed to travel to work, companies and institutions alike were encouraged to implement teleworking to avoid their employees taking public transport. After the end of the first lockdown, even though travel for personal reasons became allowed again, the teleworking measures remained in place to a certain extent and

kept on affecting public transport ridership. This is what we want to measure with our model.

Another source of information is the regression coefficients for each day of the week (fig. 3). Both segment one (corresponding to normal behavior) and segment two (during and after the pandemic) have an increased number of entries during the weekdays and a decreased ridership on weekends. But when we look at the difference between the coefficients of two segments, we notice that in several cases, the coefficients associated with weekdays decreased with the pandemic and the weekend's coefficients increased. This means that the difference between weekdays and weekends is less noticeable after the pandemic. We hypothesize that this effect is due to the use of teleworking policies that reduced the number of commuters using public transport every day. Please note that fig. 3 is not affected by modal shift, as the effect of the sum of the seven covariates is absorbed by the mean ( $\mathbf{m}_{k,s}$  in eq. 5) and not the variables encoding for the day of the week. The influence of the covariates varies from cluster to cluster, which will be helpful when understanding what makes clusters specific.

To know which features distinguish one cluster from another, the first method is to look at the regression coefficients (fig. 3). We perform a t-test on the difference of regression coefficients across a couple of clusters to know which are the coefficients that change significantly between two clusters. To compare clusters, we focus on the first segment, corresponding to a more usual behavior. Although the entries of all clusters increase during weekdays and decrease during weekends, the amplitude of weekly variations differs between clusters which is informative to us.

The first cluster has the highest variations between weekdays and weekends, indicating a higher rate of commuters than the other clusters. It includes stations from the North, West, and South suburbs, as well as the center and West Paris boroughs (results not shown), two regions with moderate to high rates of employments-to-surface [9]. As we expected, the stations with the highest weekly variations are also those where this variation dampened the most after the pandemic: the stations with the most commute are also the ones where teleworking policies are the most impactful on public transport.

The next three clusters have average weekly variations.

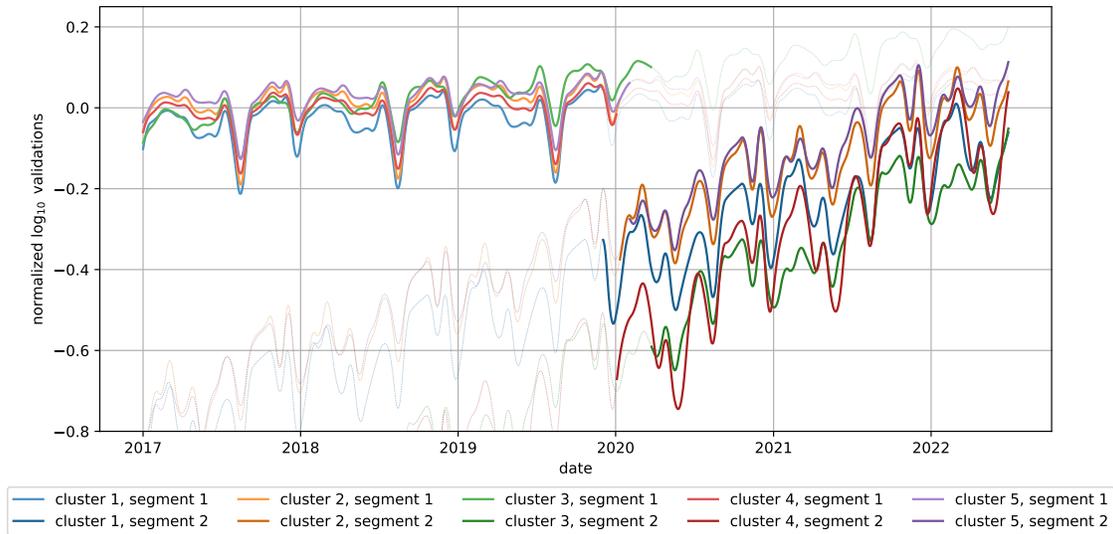


Fig. 2: The splines associated with each cluster, along with the effect of time. Dashed lines denote the continuation of each segment in regions where it is unlikely.

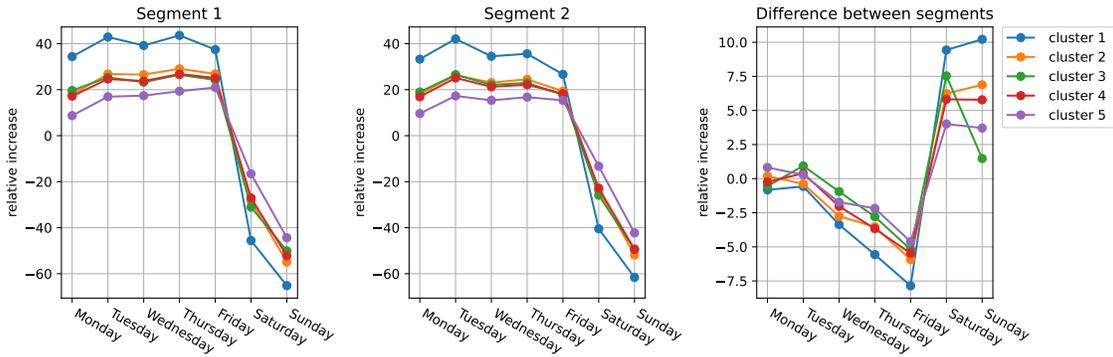


Fig. 3: The effect of the seven days of the week on each of the two segments (left, center) and the difference in regression coefficients for each cluster (all values are translated from the log scale to a linear scale). We deal with the co-linearity by computing the coefficients such that the effect of the sum of the coefficients is absorbed by each segment’s mean (which does not appear on this figure).

Cluster two is, by all accounts, close to the average, except for its variance (see table IV): both its first and second segment have particularly low variances. This means that this cluster is dedicated to stations with little noise and average weekly variations.

Conversely, cluster three has a noticeably high variance before the pandemic. This means that this cluster comprises all the stations where variations in the number of entries cannot be explained using the covariates. Surprisingly, variance in the second segment is moderate.

Cluster four stands out for its high variance after the pandemic. It groups all stations for which the number of entries cannot be accurately explained using the covariates we provided (such as the official lockdowns, for instance). It should also be noted that the decrease in the number of entries in these stations is the largest out of all five clusters, along with cluster three (fig. 2). This is compensated by the

highest recovery rate of all clusters: if we consider the impact of time on the logarithm of the number of entries, the number of entries in the second segment increases by 74% per year on average. In comparison, the second segment of cluster three increases by 60% per year, while the other three clusters increase by 37% to 42% per year during the period following the pandemic. Thus, cluster four comprising all the stations whose series of entries were affected by the pandemic.

Finally, cluster five has the smallest differences between weekdays and weekends of all the clusters. Like cluster two, it has a small variance, indicating that the stations in this cluster have variations which the model explains correctly using the exogenous variables. The decrease in the number of entries (fig. 2) is among the lowest of all clusters, along with cluster two. Geographically, the stations in this cluster are located almost exclusively in inner Paris (table IV), where housings are the most common [9].

Cluster	Variance of segment one	Variance of segment two	Average distance to the center of Paris (km)
1	$1.8 \times 10^{-2}$	$5.4 \times 10^{-2}$	16.5
2	$5.7 \times 10^{-3}$	$2.2 \times 10^{-2}$	7.2
3	$3.5 \times 10^{-2}$	$5.1 \times 10^{-2}$	16.0
4	$7.9 \times 10^{-3}$	$1.7 \times 10^{-1}$	11.8
5	$4.8 \times 10^{-3}$	$1.7 \times 10^{-2}$	5.8

TABLE IV: The variance in each segment of each cluster, along with the average distance between the stations in the cluster and the Center of Paris.

To sum up, the main element that distinguishes the clusters is the degree of variation between weekdays and weekends: clusters one and five have the largest and smallest variations, respectively. Among the middle three clusters, one of them (cluster two) had a small variance around an average behavior, while the remaining two had a high variance either before or after the pandemic. Surprisingly, no cluster had particularly high variance for both its segments. Additionally, one could have expected the consequences of the pandemic to be more diverse on the stations with high a commute, but the model groups all these stations into a single cluster.

## V. CONCLUSION

Clustering, segmentation, and regression are three well-explored research problems, to the point the intersection between any two of them has been covered in the literature. We extended the clustering-segmentation model from [18], and completed it by adding the contribution of exogenous variables. We developed the estimation of parameters using the Expectation-Maximization algorithm, and experimented with it on both synthetic and real data, demonstrating the interest in a joint estimation of parameters as opposed to a sequential one.

Using this model on a dataset of entries in the Paris public transport network allowed us to understand how the Covid-19 pandemic affected the public transport ridership. As expected, all stations saw the number of entries decrease sharply at the start of the first lockdown, with a slow recovery during the two years that followed. In addition, most clusters saw the number of entries further decrease during weekdays and increase during weekends, compared to their respective pre-pandemic levels. We attribute this effect to teleworking policies all the more so because the cluster which is defined by the highest proportion of commuters saw the strongest decrease in ridership during weekdays and the strongest increase during weekends. Whenever the variables we provided are not enough for the model, the model assigns a high variance to the concerned segment, as this was the case for another of the five clusters the model uncovered: the second segment of cluster four had a high variance, indicating that the evolution of the time series assigned to it is more irregular than the simple effects we included in with the covariates.

The next step would be to find a way to reduce the computational load of the proposed model. One could, for instance, think about initializing the EM algorithm with coefficients

found from a simpler version of the model. Another possibility for improvement would be to use model selection criteria to choose the number of segments automatically, such as penalized criteria (BIC, AIC, *etc.*) or dedicated heuristics [1].

## ACKNOWLEDGMENTS

This work was supported by Vinci through the funding of the Eco Conception Chair in École des Ponts ParisTech.

## REFERENCES

- [1] J.-P. Baudry, C. Maugis, and B. Michel. Slope heuristics: overview and implementation. *Statistics and Computing*, 22(2):455–470, Mar. 2012.
- [2] J. L. Bernal, S. Cummins, and A. Gasparrini. Interrupted time series regression for the evaluation of public health interventions: a tutorial. *International Journal of Epidemiology*, 46(1):348–355, Feb. 2017.
- [3] F. Chamroukhi. Piecewise Regression Mixture for Simultaneous Functional Data Clustering and Optimal Segmentation. *Journal of Classification*, 33(3):374–411, Oct. 2016.
- [4] F. Chang, S. Yasmin, H. Huang, A. H. S. Chan, and M. M. Haque. Injury severity analysis of motorcycle crashes: A comparison of latent class clustering and latent segmentation based models with unobserved heterogeneity. *Analytic Methods in Accident Research*, 32:100188, Dec. 2021.
- [5] A. R. Conn, N. I. M. Gould, and P. L. Toint. *Trust Region Methods*. MOS-SIAM Series on Optimization. Society for Industrial and Applied Mathematics, Jan. 2000.
- [6] P. de Nailly, E. Côme, A. Samé, L. Oukhellou, J. Ferriere, and Y. Merad-Boudia. What can we learn from 9 years of ticketing data at a major transport hub? A structural time series decomposition. *Transportmetrica A: Transport Science*, 0(0):1–25, July 2021.
- [7] R. D. De Veaux. Mixtures of linear regressions. *Computational Statistics & Data Analysis*, 8(3):227–245, Nov. 1989.
- [8] A. P. Dempster, N. M. Laird, and D. B. Rubin. Maximum Likelihood from Incomplete Data Via the EM Algorithm. *Journal of the Royal Statistical Society: Series B (Methodological)*, 39(1):1–22, 1977.
- [9] C. Etienne and O. Latifa. Model-Based Count Series Clustering for Bike Sharing System Usage Mining: A Case Study with the Vélib’ System of Paris. *ACM Transactions on Intelligent Systems and Technology*, 5(3):39:1–39:21, July 2014.
- [10] M. Leyli-Abadi, A. Samé, L. Oukhellou, N. Cheifetz, P. Mandel, C. Féliers, and O. Chesneau. Mixture of Non-homogeneous Hidden Markov Models for Clustering and Prediction of Water Consumption Time Series. In *2018 International Joint Conference on Neural Networks (IJCNN)*, pages 1–8. IEEE, 2018.
- [11] J. Liu, S. Wu, and J. V. Zidek. On Segmented Multivariate Regression. *Statistica Sinica*, 7(2):497–525, 1997. Publisher: Institute of Statistical Science, Academia Sinica.
- [12] V. M. R. Muggeo. Estimating regression models with unknown breakpoints. *Statistics in Medicine*, 22(19):3055–3071, 2003.
- [13] B. Nistal-Nuño. Segmented regression analysis of interrupted time series data to assess outcomes of a South American road traffic alcohol policy change. *Public Health*, 150:51–59, Sept. 2017.
- [14] P. Nystrup, E. Lindström, and H. Madsen. Learning hidden Markov models with persistent states by penalizing jumps. *Expert Systems with Applications*, 150:113307, July 2020.
- [15] Y. Park, Y. Choi, K. Kim, and J. K. Yoo. Machine learning approach for study on subway passenger flow. *Scientific Reports*, 12(1):2754, Feb. 2022. Number: 1 Publisher: Nature Publishing Group.
- [16] F. Picard, S. Robin, E. Lebarbier, and J.-J. Daudin. A Segmentation/Clustering Model for the Analysis of Array CGH Data. *Biometrics*, 63(3):758–766, 2007.
- [17] C. E. Rasmussen. Gaussian Processes in Machine Learning. In O. Bousquet, U. von Luxburg, and G. Rätsch, editors, *Advanced Lectures on Machine Learning: ML Summer Schools 2003, Canberra, Australia, February 2 - 14, 2003, Tübingen, Germany, August 4 - 16, 2003, Revised Lectures*, Lecture Notes in Computer Science, pages 63–71. Springer, Berlin, Heidelberg, 2004.
- [18] A. Samé, F. Chamroukhi, G. Govaert, and P. Aknin. Model-based clustering and segmentation of time series with changes in regime. *Advances in Data Analysis and Classification*, 5(4):301–321, Dec. 2011.
- [19] C. Truong, L. Oudre, and N. Vayatis. Selective review of offline change point detection methods. *Signal Processing*, 167:107299, Feb. 2020.