# Model-based clustering and segmentation of time series with changes in regime

**Allou Samé · Faicel Chamroukhi ·
Gérard Govaert · Patrice Aknin**

**Abstract** Mixture model-based clustering, usually applied to multidimensional data, has become a popular approach in many data analysis problems, both for its good statistical properties and for the simplicity of implementation of the Expectation-Maximization (EM) algorithm. Within the context of a railway application, this paper introduces a novel mixture model for dealing with time series that are subject to changes in regime. The proposed approach consists in modeling each cluster by a regression model in which the polynomial coefficients vary according to a discrete hidden process. In particular, this approach makes use of logistic functions to model the (smooth or abrupt) transitions between regimes. The model parameters are estimated by the maximum likelihood method solved by an Expectation-Maximization algorithm. The proposed approach can also be regarded as a clustering approach which operates by finding groups of time series having common changes in regime. In addition to providing a time series partition, it therefore provides a time series segmentation. The problem of selecting the optimal numbers of clusters and segments is solved by means of the Bayesian Information Criterion (BIC). The proposed approach is shown to be efficient using a variety of simulated time series and real-world time series of electrical power consumption from rail switching operations.

Allou Samé, Faicel Chamroukhi, Patrice Aknin
Institut Français des Sciences et Technologies des Transports,
de l'Aménagement et des Réseaux (IFSTTAR)
2 rue de la Butte Verte, 93166 Noisy-le-grand Cedex
E-mail: same@inrets.fr

Gérard Govaert
Université de Technologie de Compiègne (UTC)
UMR CNRS 6599, Heudiasyc
Centre de Recherches de Royallieu, BP 20529, F-60205 Compiègne Cedex
E-mail: gerard.govaert@utc.fr

# 1 Introduction

The application which gave rise to this study is an application for diagnosing problems in rail switches, that is to say the mechanisms which enable trains to change tracks at junctions. One preliminary task in the diagnostic process is identifying groups of switching operations that have similar characteristics, and this is accomplished by performing clustering on the time series of electrical power consumption, acquired during various switching operations. This kind of data is referred to in other contexts as longitudinal data [4], signals, or curves [8].

The approach adopted in this paper is mixture model-based clustering [1, 2], which has successfully been applied in numerous domains [13], and which provides, by means of the Expectation-Maximization algorithm [6], an efficient implementation framework. Typical extensions of mixture models for time series include regression mixture models [8] and random effect regression mixture models [11,7,14,12]. These approaches are based on a projection of the original time series into a space with fewer dimensions, defined by polynomial or spline basis functions. Other approaches that combine Autoregressive Moving Average (ARMA) methods and the Expectation-Maximization algorithm have also been proposed [19]. Although these approaches can be seen as an efficient way of classifying time series, all of them use a constant dynamic within each cluster; in other words, the regressive or autoregressive coefficients of the clusters do not vary with time.

However, the time series studied in this paper are subject to various changes in regime (see figure 7) as a result of the successive mechanical movements that are involved in a switching operation. Within this particular context, a specific regression model has been proposed in [3] to deal with regime changes in time series. The model in question is a regression model in which the polynomial coefficients may vary according to a discrete hidden process, and which uses logistic functions to model the (smooth or abrupt) transitions between regimes. In this paper we extend this regression model to a finite mixture model, where each cluster is represented by its own "hidden process regression model".

This paper is organized as follows. We first present a brief review of the regression mixture model for time series clustering. Then, we detail the proposed mixture model and its parameters estimation via the Expectation-Maximization (EM) algorithm [6]. Section 5 illustrates the performances of the proposed approach using simulated examples and real-world time series from an application in the railway sector.

The time series to be classified takes the form of an independent random sample $(\boldsymbol{x}_1, \ldots, \boldsymbol{x}_n)$ where each series $\boldsymbol{x}_i$ consists of a vector of $m$ random real values $(x_{i1}, \ldots, x_{im})$ observed over the fixed time grid $\boldsymbol{t} = (t_1, \ldots, t_m)$, with $t_1 < t_2 < \ldots < t_n$. The unobserved clusters corresponding to $(\boldsymbol{x}_1, \ldots, \boldsymbol{x}_n)$ will be denoted as $(z_1, \ldots, z_n)$, where $z_i \in \{1, \ldots, K\}$.

## 2 Regression mixture model for time series clustering

This section briefly recalls the regression mixture model, as formulated by Gafney and Smith [8], in the context of times series clustering.

### 2.1 Definition of the regression mixture model

Unlike standard vector-based mixture models, the density of each component of the regression mixture is represented by a polynomial "mean series" (or mean curve) parameterized by a vector of regression coefficients and a noise variance.

The regression mixture model therefore assumes that each series $\boldsymbol{x}_i$ is distributed according to the conditional mixture density

$$f(\boldsymbol{x}_i|\boldsymbol{t};\boldsymbol{\theta}) = \sum_{k=1}^{K} \pi_k \, \mathcal{N}(\boldsymbol{x}_i; \mathbf{T}\boldsymbol{\beta}_k, \sigma_k^2\mathbf{I}), \tag{1}$$

where $\boldsymbol{\theta} = (\pi_1, \ldots, \pi_K, \boldsymbol{\beta}_1, \ldots, \boldsymbol{\beta}_K, \sigma_1^2, \ldots, \sigma_K^2)$ is the complete parameter vector, the $\pi_k$ are the proportions of the mixture satisfying $\sum_{k=1}^{K} \pi_k = 1$, $\boldsymbol{\beta}_k$ and $\sigma_k^2$ are respectively the the $(p+1)$-dimensional coefficient vector of the $k$th regression model and the associated noise variance. The matrix $\mathbf{T} = (T_{uj})$ is a $m \times (p+1)$ Vandermonde matrix verifying $T_{uj} = t_j^{u-1}$ for all $1 \leq j \leq m$ and $1 \leq u \leq (p+1)$, and $\mathcal{N}(\cdot; \boldsymbol{\mu}, \boldsymbol{\Sigma})$ is the Gaussian density with mean vector $\boldsymbol{\mu}$ and covariance matrix $\boldsymbol{\Sigma}$.

### 2.2 Fitting the model

Assuming that the sample $(\boldsymbol{x}_1, \ldots, \boldsymbol{x}_n)$ is independent, the parameter vector $\boldsymbol{\theta}$ is estimated by maximizing the conditional log-likelihood

$$\mathcal{L}(\boldsymbol{\theta}) = \log \sum_{i=1}^{n} f(\boldsymbol{x}_i|\boldsymbol{t};\boldsymbol{\theta}) = \sum_{i=1}^{n} \log \sum_{k=1}^{K} \pi_k \mathcal{N}(\boldsymbol{x}_i; \mathbf{T}\boldsymbol{\beta}_k, \sigma_k^2\mathbf{I}) \tag{2}$$

via the Expectation-Maximization (EM) algorithm initiated by Dempster, Laird and Rubin [6].

Once the parameters have been estimated, a time series partition is obtained by assigning each series $\boldsymbol{x}_i$ to the cluster having the highest posterior probability

$$p(z_i = k|\boldsymbol{t}, \boldsymbol{x}_i; \boldsymbol{\theta}) = \frac{\pi_k \mathcal{N}(\boldsymbol{x}_i; \mathbf{T}\boldsymbol{\beta}_k, \sigma_k^2\mathbf{I})}{\sum_{h=1}^{K} \pi_h \mathcal{N}(\boldsymbol{x}_i; \mathbf{T}\boldsymbol{\beta}_h, \sigma_h^2\mathbf{I})}. \tag{3}$$

## 3 Clustering time series with changes in regime

3.1 The global mixture model

As with the standard regression mixture model, the mixture model introduced for clustering time series with changes in regime assumes that the series $\boldsymbol{x}_i$ are independently generated according to the global mixture model

$$f(\boldsymbol{x}_i|\boldsymbol{t};\boldsymbol{\theta}) = \sum_{k=1}^{K} \pi_k f_k(\boldsymbol{x}_i|\boldsymbol{t};\boldsymbol{\theta}_k), \tag{4}$$

where $\boldsymbol{\theta} = (\pi_1,\dots,\pi_K,\boldsymbol{\theta}_1,\dots,\boldsymbol{\theta}_K)$, $\pi_1,\dots,\pi_K$ denote the proportions of the mixture, and $\boldsymbol{\theta}_k$ the parameters of the different component densities $f_k$. The main difference between the model proposed here and Gafney and Smith's regression mixture model [8] lies in the definition of the component densities $f_k$, described in the following section.

3.2 Definition of the mixture components

We assume that the $k$th cluster, that is to say the time series corresponding to the component $f_k$ of the proposed mixture, is generated as follows. Given the cluster label $z_i = k$ and the fixed time vector $\boldsymbol{t}$, a time series $\boldsymbol{x}_i$ is generated according to a specific regression model which implicitly supposes that there are $L$ $p$th order polynomial regression models involved in the generation of $\boldsymbol{x}_i$. The assignment of the $x_{ij}$'s to the different (sub) regression models is specified by a hidden process denoted by $\boldsymbol{w}_i = (w_{i1},\dots,w_{im})$, where $w_{ij} \in \{1,\dots,L\}$. Thus, given the cluster label $z_i = k$, the individual observations $x_{ij}$ of a series $\boldsymbol{x}_i$ are generated as follows:

$$\forall j = 1,\dots,m, \quad \begin{cases} x_{ij} = \sum_{\ell=1}^{L} w_{ij\ell}\big(\mathbf{T}_j'\boldsymbol{\beta}_{k\ell} + \sigma_{k\ell}\varepsilon_{ij}\big) \\ \varepsilon_{ij} \sim \mathcal{N}(0,1) \end{cases}, \tag{5}$$

where $w_{ij\ell} = 1$ if $w_{ij} = \ell$ and 0 otherwise. The parameters $\sigma_{k\ell}$ and $\boldsymbol{\beta}_{k\ell}$ are respectively the noise standard deviation and the $(p+1)$-dimensional coefficient vector of the $\ell$th regression model of the $k$th cluster. $\mathbf{T}_j'$ denotes the transpose of the vector $\mathbf{T}_j = (1, t_j, \dots, t_j^p)^T$.

The regression component labels $w_{ij}$ $(j = 1,\dots,m)$ are assumed to be generated according to the multinomial distribution $\mathcal{M}(1, \pi_{k1}(t_j;\boldsymbol{\alpha}_k),\dots,\pi_{kL}(t_j;\boldsymbol{\alpha}_k))$, where

$$\pi_{k\ell}(t;\boldsymbol{\alpha}_k) = \frac{\exp(\boldsymbol{\alpha}_{k\ell1}t + \boldsymbol{\alpha}_{k\ell0})}{\sum_{h=1}^{L} \exp(\boldsymbol{\alpha}_{kh1}t + \boldsymbol{\alpha}_{kh0})}. \tag{6}$$

is a logistic function with parameter vector $\boldsymbol{\alpha}_k = \{\boldsymbol{\alpha}_{k\ell}; \ell = 1,\dots,L\}$ and $\boldsymbol{\alpha}_{k\ell} = (\boldsymbol{\alpha}_{k\ell0}, \boldsymbol{\alpha}_{k\ell1})$. A logistic function defined in this way ensures a smooth transition between the different polynomial regimes. Thus, given $z_i = k$ and

$t_j$, the individual observations $x_{ij}$ of a series $\boldsymbol{x}_i$ are independently distributed according to the mixture model given by

$$p(x_{ij}|t_j; \boldsymbol{\theta}_k) = \sum_{\ell=1}^{L} \pi_{k\ell}(t_j; \boldsymbol{\alpha}_k)\mathcal{N}(x_{ij}; \boldsymbol{\beta}_{k\ell}^T\mathbf{T}_j, \sigma_{k\ell}^2). \qquad (7)$$
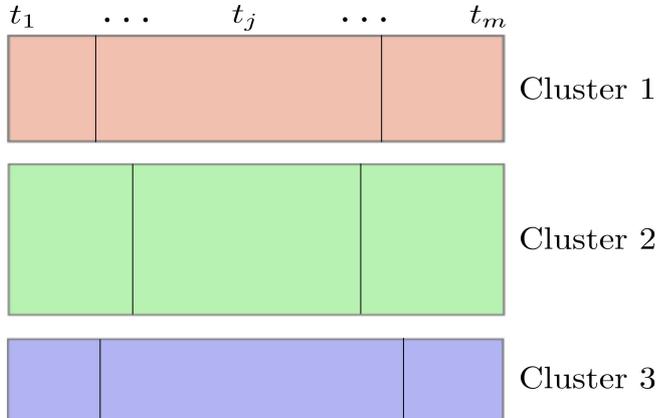
The density $f_k$ can thus be written as

$$f_k(\boldsymbol{x}_i|\boldsymbol{t}; \boldsymbol{\theta}_k) = \prod_{j=1}^{m}\sum_{\ell=1}^{L} \pi_{k\ell}(t_j; \boldsymbol{\alpha}_k)\mathcal{N}(x_{ij}; \boldsymbol{\beta}_{k\ell}^T\mathbf{T}_j, \sigma_{k\ell}^2). \qquad (8)$$

3.3 A cluster-segmentation model

The proposed model leads to the segmentation $\mathbf{E}_k = (E_{k\ell})_{\ell=1,\ldots,L}$ of the set of time series originating from the $k$th cluster, where

$$E_{k\ell} = \left\{ t \in [t_1; t_m] \ / \ \pi_{k\ell}(t; \boldsymbol{\alpha}_k) = \max_{1 \leq h \leq L} \pi_{kh}(t; \boldsymbol{\alpha}_k) \right\}. \qquad (9)$$

It can be proved that the set $E_{k\ell}$ is convex (see appendix A). Therefore, $\mathbf{E}_k$ is a segmentation into contiguous parts of $\{t_1, \ldots t_m\}$. Figure 1 illustrates the latent structure of the proposed model with $K = 3$ and $L = 3$.



**Fig. 1** Latent hierarchical structure of the proposed model with $K = 3$ clusters: for each time series cluster, the vertical lines define a segmentation into $L = 3$ segments

3.4 Parameter estimation via the EM algorithm

The parameters of the proposed model are estimated by maximizing the conditional log-likelihood defined by

$$
\begin{aligned}
\mathcal{L}(\boldsymbol{\theta}) &= \sum_{i=1}^{n} \log f(\boldsymbol{x}_i | \boldsymbol{t}; \boldsymbol{\theta}) \\
&= \sum_{i=1}^{n} \log \sum_{k=1}^{K} \pi_k \Big( \prod_{j=1}^{m} \sum_{\ell=1}^{L} \pi_{k\ell}(t_j; \boldsymbol{\alpha}_k) \mathcal{N}(x_{ij}; \boldsymbol{\beta}_{k\ell}^{T} \mathbf{T}_j, \sigma_{k\ell}^2) \Big). \quad (10)
\end{aligned}
$$

The Expectation Maximization (EM) algorithm [6] is used for the maximization of this log-likelihood, a problem which cannot be solved analytically. Let us recall that the EM algorithm requires a complete data specification, whose log-likelihood can be maximized more easily than the observed data log-likelihood. Here, the "complete data" are obtained by adding to each series $\boldsymbol{x}_i$ its cluster membership $z_i$ and its assignment process $\boldsymbol{w}_i = (w_{ij})_{j=1,\dots,m}$ to the different sub-regression models. Using the binary coding of $z_i$ and $\boldsymbol{w}_{ij}$,

$$
z_{ik} = \left\{ \begin{array}{ll} 1 & \text{if } z_i = k \\ 0 & \text{otherwise} \end{array} \right. \quad \text{and} \quad w_{ij\ell} = \left\{ \begin{array}{ll} 1 & \text{if } w_{ij} = \ell \\ 0 & \text{otherwise,} \end{array} \right.
$$

the complete data log-likelihood can be written as

$$
\begin{aligned}
\mathcal{L}_c(\boldsymbol{\theta}) &= \sum_{i=1}^{n} \log p(\boldsymbol{x}_i, z_i, \boldsymbol{w}_i | \boldsymbol{t}; \boldsymbol{\theta}) = \sum_{i=1}^{n} \sum_{k=1}^{K} z_{ik} \log \pi_k + \\
&\quad \sum_{i=1}^{n} \sum_{j=1}^{m} \sum_{k=1}^{K} \sum_{\ell=1}^{L} z_{ik} w_{ij\ell} \log \Big( \pi_{k\ell}(t_j; \boldsymbol{\alpha}_k) \mathcal{N}(x_{ij}; \boldsymbol{\beta}_{k\ell}^{T} \mathbf{T}_j, \sigma_{k\ell}^2) \Big). \quad (11)
\end{aligned}
$$

Given an initial value of the parameter vector $\boldsymbol{\theta}^{(0)}$, the EM algorithm alternates the two following steps until convergence.

*E-Step (Expectation)*

This step consists in evaluating the expectation of the complete data log-likelihood conditionally on the observed data and the current parameter vector $\boldsymbol{\theta}^{(q)}$, $q$ denoting the current iteration:

$$
\begin{aligned}
Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{(q)}) &= E \Big[ \mathcal{L}_c(\theta) \big| \boldsymbol{t}, \boldsymbol{x}_1, \dots, \boldsymbol{x}_n; \boldsymbol{\theta}^{(q)} \Big] = \sum_{i=1}^{n} \sum_{k=1}^{K} r_{ik}^{(q)} \log \pi_k + \\
&\quad \sum_{i=1}^{n} \sum_{j=1}^{m} \sum_{k=1}^{K} \sum_{\ell=1}^{L} \lambda_{ijk\ell}^{(q)} \log \Big( \pi_{k\ell}(t_j; \boldsymbol{\alpha}_k) \mathcal{N}(x_{ij}; \boldsymbol{\beta}_{k\ell}^{T} \mathbf{T}_j, \sigma_{k\ell}^2) \Big) \quad (12)
\end{aligned}
$$

where

$$r_{ik}^{(q)} = E[z_{ik}|\boldsymbol{t}, \boldsymbol{x}_i; \boldsymbol{\theta}^{(q)}] = \frac{\pi_k^{(q)} f_k(\boldsymbol{x}_i|\boldsymbol{t}; \boldsymbol{\theta}_k^{(q)})}{\sum_{h=1}^{K} \pi_h^{(q)} f_h(\boldsymbol{x}_i|\boldsymbol{t}; \boldsymbol{\theta}_h^{(q)})} \tag{13}$$

is the posterior probability that time series $\boldsymbol{x}_i$ originates from cluster $k$, and

$$\begin{aligned}
\lambda_{ijk\ell}^{(q)} &= E[z_{ik}\, w_{ij\ell}|\boldsymbol{t}, \boldsymbol{x}_i; \boldsymbol{\theta}^{(q)}] \\
&= \frac{\pi_k^{(q)} f_k(\boldsymbol{x}_i|\boldsymbol{t}; \boldsymbol{\theta}_k^{(q)})}{\sum_{h=1}^{K} \pi_h^{(q)} f_h(\boldsymbol{x}_i|\boldsymbol{t}; \boldsymbol{\theta}_h^{(q)})} \times \frac{\pi_{k\ell}(t_j; \boldsymbol{\alpha}_k^{(q)}) \mathcal{N}(x_{ij}; \boldsymbol{\beta}_{k\ell}^{(q)^T} \mathbf{T}_j, \sigma_{k\ell}^{2(q)})}{\sum_{h=1}^{L} \pi_{kh}(t_j; \boldsymbol{\alpha}_k^{(q)}) \mathcal{N}(x_{ij}; \boldsymbol{\beta}_{kh}^{(q)^T} \mathbf{T}_j, \sigma_{kh}^{2(q)})}
\end{aligned} \tag{14}$$

is the posterior probability that $(t_j, x_{ij})$ originates from the $\ell$th sub-regression model of cluster $k$.

*M-Step (Maximization)*

This step consists in computing the parameter vector $\boldsymbol{\theta}^{(q+1)}$ that maximizes the quantity $Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{(q)})$ with respect to $\boldsymbol{\theta}$. For our purposes this quantity can be written as

$$Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{(q)}) = Q_1((\pi_k)) + Q_2((\boldsymbol{\alpha}_k)) + Q_3((\beta_{k\ell}, \sigma_{k\ell}^2)),$$

where

$$Q_1((\pi_k)) = \sum_{i=1}^{n} \sum_{k=1}^{K} r_{ik}^{(q)} \log \pi_k, \tag{15}$$

$$Q_2((\boldsymbol{\alpha}_k)) = \sum_{i=1}^{n} \sum_{j=1}^{m} \sum_{k=1}^{K} \sum_{\ell=1}^{L} \lambda_{ijk\ell}^{(q)} \log \left(\pi_{k\ell}(t_j; \boldsymbol{\alpha}_k)\right), \tag{16}$$

$$Q_3((\beta_{k\ell}, \sigma_{k\ell}^2)) = \sum_{i=1}^{n} \sum_{j=1}^{m} \sum_{k=1}^{K} \sum_{\ell=1}^{L} \lambda_{ijk\ell}^{(q)} \log \left(\mathcal{N}(x_{ij}; \boldsymbol{\beta}_{k\ell}^T \mathbf{T}_j, \sigma_{k\ell}^2)\right). \tag{17}$$

$Q$ can thus be maximized by separately maximizing the quantities $Q_1$, $Q_2$ and $Q_3$. As in the classical Gaussian mixture model, it can easily be shown that the proportions $\pi_k$ that maximize $Q_1$ under the constraint $\sum_{k=1}^{K} \pi_k = 1$ are given by

$$\pi_k^{(q+1)} = \frac{\sum_{i=1}^{n} r_{ik}^{(q)}}{n}. \tag{18}$$

$Q_2$ can be maximized with respect to the $\boldsymbol{\alpha}_k$ by separately solving $K$ weighted logistic regression problems:

$$\boldsymbol{\alpha}_k^{(q+1)} = \arg\max_{\boldsymbol{\alpha}_k} \sum_{i=1}^{n} \sum_{j=1}^{m} \sum_{\ell=1}^{L} \lambda_{ijk\ell}^{(q)} \log \left(\pi_{k\ell}(t_j; \boldsymbol{\alpha}_k)\right) \tag{19}$$

through the well known Iteratively Reweighted Least Squares (IRLS) algorithm [9,3]. Let us recall that the IRLS algorithm, which is generally used

to estimate the parameters of a logistic regression model, is equivalent to the following Newton Raphson algorithm [9,3]:

$$\boldsymbol{\alpha}_k^{(v+1)} = \boldsymbol{\alpha}_k^{(v)} - \Big[\frac{\partial^2 Q_{2k}}{\partial \boldsymbol{\alpha}_k \partial \boldsymbol{\alpha}_k^T}\Big]^{-1}_{\boldsymbol{\alpha}_k = \boldsymbol{\alpha}_k^{(v)}} \Big[\frac{\partial Q_{2k}}{\partial \boldsymbol{\alpha}_k}\Big]_{\boldsymbol{\alpha}_k = \boldsymbol{\alpha}_k^{(v)}}, \tag{20}$$

where

$$Q_{2k} = \sum_{i=1}^{n} \sum_{j=1}^{m} \sum_{\ell=1}^{L} \lambda_{ijk\ell}^{(q)} \log \pi_{k\ell}(t_j; \boldsymbol{\alpha}_k).$$

Maximizing $Q_3$ with respect to $\boldsymbol{\beta}_{k\ell}$ consists in analytically solving $K \times L$ weighted least-squares problems. It can be shown that

$$\boldsymbol{\beta}_{k\ell}^{(q+1)} = \Big[\mathbf{T}'\big(\sum_{i=1}^{n} \Lambda_{ik\ell}^{(q)}\big)\mathbf{T}\Big]^{-1} \Big[\mathbf{T}\big(\sum_{i=1}^{n} \Lambda_{ik\ell}^{(q)}\big)\boldsymbol{x}_i\Big], \tag{21}$$

where $\Lambda_{ik\ell}^{(q)}$ is the $m \times m$ diagonal matrix whose diagonal elements are $\{\lambda_{ijk\ell}^{(q)} \; ; \; j = 1, \ldots, m\}$. The maximization of $Q_3$ with respect to $\sigma_{k\ell}^2$ gives

$$(\sigma_{k\ell}^2)^{(q+1)} = \frac{\sum_{i=1}^{n} \big\|\sqrt{\Lambda_{ik\ell}^{(q)}}(\boldsymbol{x}_i - \mathbf{T}\boldsymbol{\beta}_{k\ell}^{(q+1)})\big\|^2}{\sum_{i=1}^{n} \operatorname{trace}(\Lambda_{ik\ell}^{(q)})}, \tag{22}$$

where $\sqrt{\Lambda_{ik\ell}^{(q)}}$ is the $m \times m$ diagonal matrix whose diagonal elements are $\{\sqrt{\lambda_{ijk\ell}^{(q)}} \; ; \; j = 1, \ldots, m\}$ and $\| \cdot \|$ is the norm corresponding to the euclidian distance.

*M-step for three parsimonious models*

**Common segmentation for all clusters**
In certain situations, the segmentation defined by the $\boldsymbol{\alpha}_k$ $(k = 1, \ldots, K)$ may be constrained to be common for each cluster, that is $\boldsymbol{\alpha}_k = \boldsymbol{\alpha} \; \forall k$. In that case, the quantity $Q_2$ to be maximized can be rewritten as:

$$Q_2(\boldsymbol{\alpha}) = \sum_{i=1}^{n} \sum_{j=1}^{m} \sum_{\ell=1}^{L} \lambda_{ij \cdot \ell}^{(q)} \log\Big(\pi_\ell(t_j; \boldsymbol{\alpha})\Big), \tag{23}$$

where $\lambda_{ij \cdot \ell}^{(q)} = \sum_{k=1}^{K} \lambda_{ijk\ell}^{(q)}$. The IRLS algorithm can therefore be used to compute the parameter $\boldsymbol{\alpha}^{(q+1)}$, in the same way as for the unconstrained situation.

**Common variance for regression models from the same cluster**

In other situations, it may be useful to constrain the regression models variances to be common within a same cluster. In that case, $\sigma_{k\ell}^2 = \sigma_k^2 \ \forall k, \ell$. The updating formula for the variance can thus be written as:

$$\left(\sigma_k^2\right)^{(q+1)} = \frac{\sum_{i=1}^n \sum_{\ell=1}^L \left\| \sqrt{\Lambda_{ik\ell}^{(q)}} \left(\boldsymbol{x}_i - \mathbf{T}\boldsymbol{\beta}_{k\ell}^{(q+1)}\right) \right\|^2}{\sum_{i=1}^n \sum_{\ell=1}^L \text{trace}(\Lambda_{ik\ell}^{(q)})}. \tag{24}$$

**Common variance for all regression models**

If the model variances are constrained to be common all regression models, we have $\sigma_{k\ell}^2 = \sigma^2 \ \forall k, \ell$. The updating formula for the variance takes the form:

$$\left(\sigma^2\right)^{(q+1)} = \frac{\sum_{i=1}^n \sum_{k=1}^K \sum_{\ell=1}^L \left\| \sqrt{\Lambda_{ik\ell}^{(q)}} \left(\boldsymbol{x}_i - \mathbf{T}\boldsymbol{\beta}_{k\ell}^{(q+1)}\right) \right\|^2}{n \times m}. \tag{25}$$

3.5 Time series clustering, approximation and segmentation

From the parameters estimated by the EM algorithm, a partition of the time series can easily be deduced by applying the maximum a posteriori (MAP) rule

$$z_i = \arg\max_k r_{ik}. \tag{26}$$

The clusters "mean series" can be approximated by the series $\mathbf{c}_k = (c_{kj})$, with

$$c_{kj} = E[x_{ij}|t_j, z_i = k; \boldsymbol{\theta}] = \sum_{\ell=1}^L \pi_{k\ell}(t_j; \boldsymbol{\alpha}_k)\mathbf{T}_j'\boldsymbol{\beta}_{k\ell}. \tag{27}$$

Moreover, a segmentation $\mathbf{E}_k = (E_{k\ell})_{\ell=1,\dots,L}$ of the time series originating from the $k$th cluster can be derived from the estimated parameters by computing $E_{k\ell}$ as defined in equation 9.

3.6 Assessing the number of clusters, segments and the regression order

In the context of mixture models and the EM algorithm, the natural criterion for model selection is the Bayesian Information Criterion (BIC) [16]. Unlike for classical mixture regression models, three parameters need to be tuned: the number of clusters $K$, the number of segments $L$ and the degree $p$ of the polynomials. The BIC criterion, in this case, can be defined by:

$$BIC(K, L, p) = L(\widehat{\boldsymbol{\theta}}) - \frac{\nu(K, L, p)}{2}\log(n), \tag{28}$$

where $\widehat{\boldsymbol{\theta}}$ is the parameter vector estimated by the EM algorithm, and $\nu(K, L, p)$ is the number of free parameters of the model. In the proposed model, the number of free parameters

$$\nu(K, L, p) = (K - 1) + 2\,K(L - 1) + L\,K(p + 1) + L\,K \qquad (29)$$

is the sum of the mixture proportions, the logistic functions parameters, the polynomial coefficients and the variances.

From a practical point of view, the maximum numbers $K_{max}$, $L_{max}$ and $p_{max}$ are first specified. Then, the EM algorithm is run for $K \in \{1, \ldots, K_{max}\}$, $L \in \{1, \ldots, L_{max}\}$ and $p \in \{1, \ldots, p_{max}\}$, and the BIC criterion is computed. The set $(K, L, p)$ with the highest value of BIC is taken to be right solution.

## 4 Experimental study

This section is devoted to an evaluation of the clustering accuracy of the proposed algorithm, carried out using simulated time series and real-world time series from a railway application. Results obtained from the proposed algorithm are compared with those provided by the clustering approach, based on the regression mixture described in section 2. To measure the clustering accuracy, two criteria were used: the misclassification percentage between the true partition and the estimated partition, and the intra-cluster inertia $\sum_{k=1}^{K} \sum_{i=1}^{n} \widehat{z}_{ik} ||\boldsymbol{x}_i - \widehat{\mathbf{c}}_k||^2$, where $(\widehat{z}_{ik})$ and $\widehat{\mathbf{c}}_k = (\widehat{c}_{kj})_{j=1,\ldots,m}$ represent respectively the binary partition matrix and the $k$th mean series estimated by each of the two compared algorithms:

- $\widehat{c}_{kj} = \sum_{\ell=1}^{L} \pi_{k\ell}(t_j; \boldsymbol{\alpha}_k)\mathbf{T}'_j\boldsymbol{\beta}_{k\ell}$ for the proposed algorithm,
- $\widehat{c}_{kj} = \mathbf{T}'_j\boldsymbol{\beta}_k$ for the regression mixture EM algorithm.

4.1 Experiments using simulated data

*4.1.1 Simulation protocol and algorithms tuning*

The time series are simulated as follows: $n$ series of length $m$ are generated according to a mixture of $K$ clusters whose mean series can be either polynomial or the sum of polynomials weighted by logistic functions.
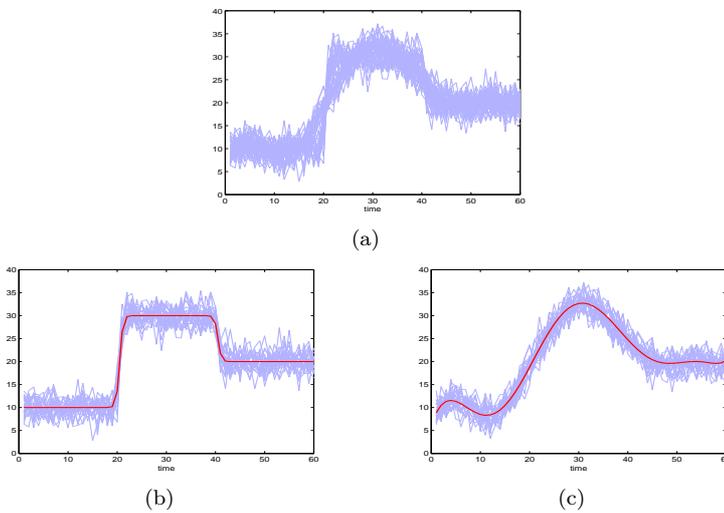
The polynomial coefficients and variances are initialized as follows: $K$ series are randomly selected and segmented into $L$ regularly spaced segments; the polynomial regression parameters are derived from a $p$th order regression on each segment. The logistic regression parameters are initialized to the null vector. The initial polynomial coefficients and variances of the regression mixture approach are obtained by performing a $p$th-order regression on $K$ randomly drawn series. The proportions of the initial clusters are set to $1/K$ for all algorithms. Each algorithm starts with 20 different initializations and the solution with the highest log-likelihood is selected.

### 4.1.2 Comparison between the proposed model and the standard regression mixture

The experiments were performed in order to compare the relative performances of the proposed EM algorithm and the EM algorithm for standard regression mixtures. So as not to favor either method unduly, the data were generated without reference either to the proposed model or to the regression mixture. Each data set, consisting of $n = 50$ time series of length $m = 60$, was simulated according to a mixture of $K = 2$ clusters with equal proportions ($\pi_1 = \pi_2 = 1/2$). The first cluster mean curve was built from three polynomials of degree $p = 0$ weighted by logistic functions, while the second was a single polynomial of degree $p = 8$. Values of the variance $\sigma_{k\ell}^2$ were chosen equal for each of the simulated sets of time series. The parameters of the mean curves are given in table 1, and figure 2 provides an illustration of time series simulated according to this model.
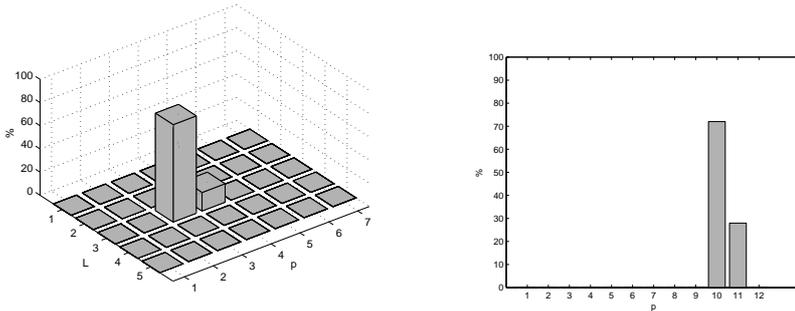
**Table 1** Clusters' mean series with their parameters

| Cluster | Mean series | Parameters | |
|---|---|---|---|
| $k = 1$ | $c_{1j} = \sum_{\ell=1}^{3} \pi_{1\ell}(t_j; \boldsymbol{\alpha}_1)\mathbf{T}_j'\boldsymbol{\beta}_{1\ell}$ | $\boldsymbol{\beta}_{11} = 10$ | $\boldsymbol{\alpha}_{11} = (1039, -34.4)'$ |
| | | $\boldsymbol{\beta}_{12} = 20$ | $\boldsymbol{\alpha}_{12} = (677, -16.7)'$ |
| | | $\boldsymbol{\beta}_{13} = 30$ | $\boldsymbol{\alpha}_{13} = (0, 0)'$ |
| $k = 2$ | $c_{2j} = \mathbf{T}_j'\boldsymbol{\beta}_2$ | $\boldsymbol{\beta}_2 = (7.4, 1.9, -0.3, -2 \times 10^{-3}, 2 \times 10^{-4},$ | |
| | | $-1.3 \times 10^{-4}, 3.2 \times 10^{-6}, -3.7 \times 10^{-8}, 1.6 \times 10^{-10})'$ | |



(a)



(b)                                    (c)

**Fig. 2** Example of $n = 50$ simulated time series (a) and series corresponding to the two clusters, with their mean (b and c)

Preliminarily, the triplet $(K, L, p)$ for the proposed approach is tuned using the BIC criterion as follows: (i) twenty-five sets of 50 time series are randomly generated with $\sigma_k^2 = 2$ ; (ii) the proposed algorithm is run on each data set, with $K \in \{1, \ldots, K_{max}\}$, $L \in \{1, \ldots, L_{max}\}$ and $p \in \{1, \ldots, p_{max}\}$ ; (iii) the selection rate for each triplet $(K, L, p)$ over the 25 random samples is computed as a percentage. The model with the highest percentage of selections is the one with $(K, L, p) = (2, 3, 3)$. The same strategy was applied to the regression mixture approach, where the pair $(K, p) = (2, 10)$ was found to have the highest percentage of selections. Figure 3 shows the percentages obtained with the two algorithms, only for $K = 2$.



**Fig. 3** Percentage of selecting respectively $(L, p)$ and $p$ by the BIC criterion for the proposed approach (left) and the regression mixture approach (right), with $K = 2$
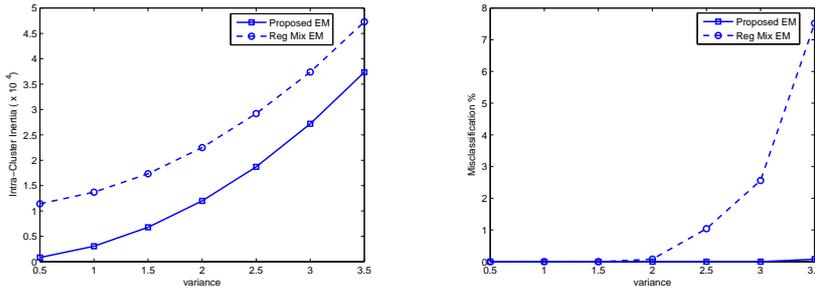
Using the optimal numbers of clusters, segments and polynomial orders computed above, the two algorithms are then compared. Table 2 gives the obtained misclassification percentages and intra-cluster inertia averaged over 25 random samples. The overall performance of the proposed algorithm is seen to be better than that of the regression mixture EM algorithm.

**Table 2** Misclassification rate and intra-cluster inertia obtained with the two compared algorithms

|                    | Misclassification percentage | Intra-cluster inertia |
| ------------------ | ---------------------------- | --------------------- |
| Proposed approach  | 0                            | $1.20 \times 10^4$    |
| Regression mixture | 0.08                         | $2.25 \times 10^4$    |

Figure 4 shows the misclassification percentage and the intra-cluster inertia (averaged over 25 different random samples of time series) in relation to the variance $\sigma_k^2$, obtained with the proposed algorithm and the regression mixture EM algorithm. The proposed algorithm is seen to outperform its competitor. Although the misclassification percentages of the two approaches are close in particular for $\sigma_k^2 \leq 2$, the intra-cluster inertia differs from about $10^4$. Misclassification provided by the regression mixture EM algorithm increases for
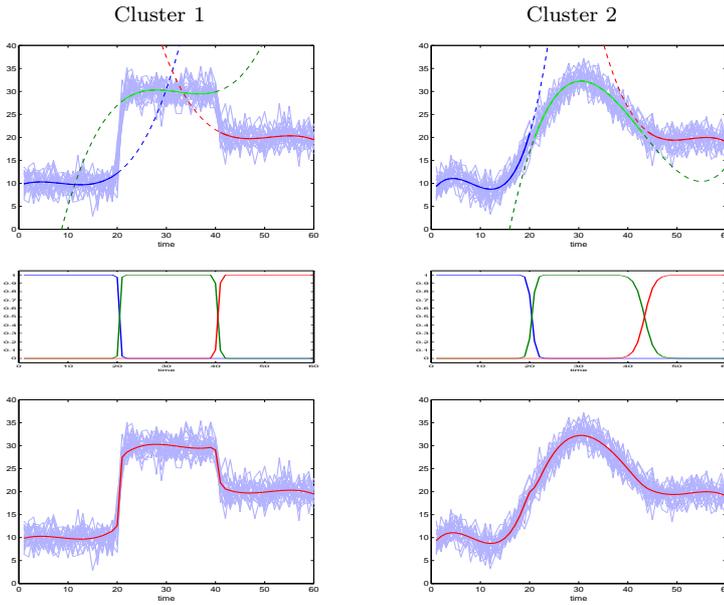
variances greater than 2.5. The intra-cluster inertia obtained by the two approaches naturally increases with the variance level, but the proposed approach performs better than its competitor. Examples of clustering results provided by the proposed approach are displayed in figure 5. It will be observed that our approach is also capable of modeling the cluster 2, whose mean series is a polynomial of degree 8, by means of three polynomials of order 3 weighted by logistic functions. Figure 6 illustrates that the regression mixture model, in contrast to the proposed model, cannot accurately model cluster 1, whose series are subject to changes in regime.
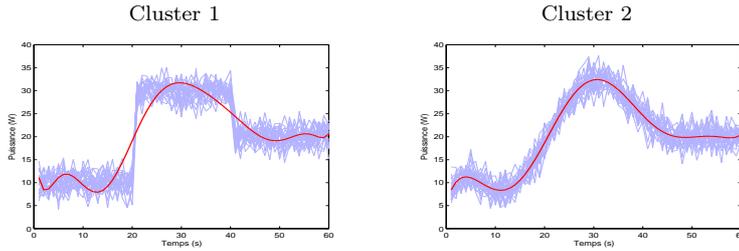


**Fig. 4** Misclassification rate (left) and intra-cluster inertia (right) in relation to the noise variance, obtained with the proposed EM algorithm and the standard regression mixture EM algorithm

4.2 Experiments using real world data

As mentioned in the introduction, the main motivation behind this study was diagnosing problems in the rail switches that allow trains to change tracks at junctions. An important preliminary task in the diagnostic process is the automatic identification of groups of switching operations that have similar characteristics, by analyzing time series of electrical power consumption acquired during switching operations. The specificity of the time series to be analyzed in this context is that they are subject to various changes in regime as a result of the mechanical movements involved in a switching operation. We accomplished this clustering task using our EM algorithm, designed for estimating the parameters of a mixture of hidden process regression models. We compared the proposed EM algorithm to the regression mixture EM algorithm previously described, on a data set of $n = 140$ time series (see figure 7). This data set is composed of four clusters identified by an expert: a defect-free cluster (35 time series), a cluster with a minor defect (40 time series), a cluster with a type 1 critical defect (45 time series) and a cluster with a type 2 critical defect (20 time series).
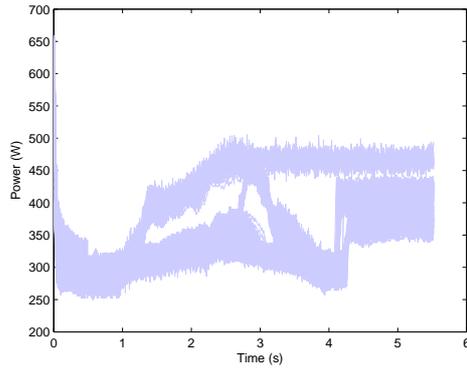
**Fig. 5** Clustering results provided by the proposed EM algorithm applied with $K = 2$, $L = 3$ and $p = 3$: clusters with their estimated polynomials (top), logistic probabilities (middle), clusters with their mean series (bottom)



**Fig. 6** Clusters and mean series estimated by the regression mixture EM algorithm applied with $(K, p) = (2, 10)$
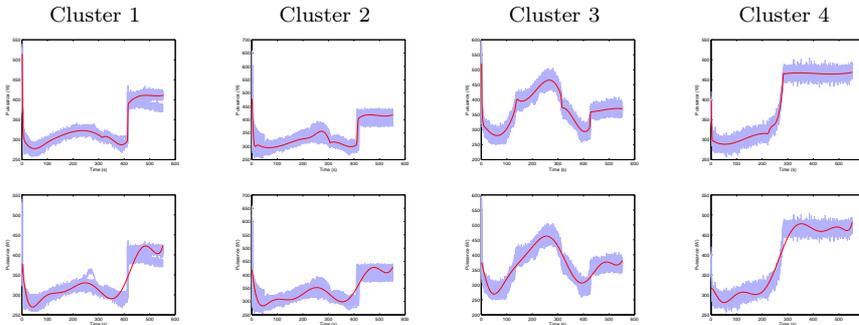
The number of regression components of the proposed algorithm was set to $L = 5$ in accordance with the number of mechanical phases in a switching operation, and the degree of the polynomial regression $p$ was set to 3, which is more appropriate for the different regimes in the time series. The polynomial order for the regression mixture approach was set to $p = 10$ which, in practice, gives the best error rates. For all the compared algorithms the number of clusters was set to $K = 4$. Table 3 shows the misclassification error rates and the corresponding intra-cluster inertia. It can be seen that the proposed regression approach provides the smallest intra-cluster error and misclassification rate. Figure 8 displays the clusters provided by the three compared algorithms and their estimated mean series.

**Fig. 7** Electrical power consumption time series acquires during $n = 140$ switch operations

**Table 3** Error obtained for the three compared approaches

|  | Regression mixture EM | Proposed EM |
|---|---|---|
| Misclassification % | 11.42 | 9.28 |
| Intra-cluster inertia | $2.6583 \times 10^{7}$ | $1.1566 \times 10^{7}$ |



**Fig. 8** Clusters and mean series estimated by the proposed EM algorithm (top) and the regression mixture EM algorithm (bottom)

## 5 Conclusion

A new mixture model-based approach for the clustering of univariate time series with changes in regime has been proposed in this paper. This approach involves modeling each cluster using a particular regression model whose polynomial coefficients vary over time according to a discrete hidden process. The transition between regimes is smoothly controlled by logistic functions. The model parameters are estimated by the maximum likelihood method, solved by a dedicated Expectation-Maximization (EM) algorithm. The proposed approach can also be regarded as a clustering approach which operates by finding groups of time series having common changes in regime. The Bayesian Information Criterion (BIC) is used to determine the numbers of clusters and

segments, as well as the regression order. The experimental results, both from simulated time series and from a real-world application, show that the proposed approach is an efficient means for clustering univariate time series with changes in regime.

## References

1. J. D. Banfield and A. E. Raftery, Model-based Gaussian and non Gaussian clustering, Biometrics, 49, 803-821, (1992)
2. G. Celeux and G. Govaert, Gaussian parsimonious clustering models, Pattern Recognition 28 (5), 781-793, (1995)
3. F. Chamroukhi, A. Samé, G. Govaert, P. Aknin, A hidden process regression model for functional data description: Application to curve discrimination, Neurocomputing, 73(7-9), 1210-1221, (2010)
4. J. M. Chiou and P. L. Li, Functional clustering and identifying sbstructures of longitudinal data, Journal of the Royal Stastical Society, Series B, 69, 679-699, (2007)
5. G. Coke and M. Tsao, Random effects mixture models for clustering electrical load series, Journal of time series analysis, 31(6), 451-464, (2010)
6. A. P. Dempster and N.M. Laird and D.B. Rubin, Maximum likelihood from incomplete data via the EM algorithm, Journal of the Royal Statistical Society, Series B, 39(1), 1-38, (1977)
7. S. J. Gaffney, P. Smyth, Curve Clustering with Random Effects Regression Mixtures, Proceedings of the Ninth International Workshop on Artificial Intelligence and Statistics, (2003)
8. S. J. Gaffney, P. Smyth, Trajectory Clustering with Mixtures of Regression Models, Proceedings of the fifth ACM SIGKDD international conference on Knowledge discovery and data mining, (1999)
9. P. Green, Iteratively Reweighted Least Squares for Maximum Likelihood Estimation, and some robust and resistant alternatives, Journal of the Royal Statistical Society, B, 46(2): 149-192, (1984)
10. G. Hébrail, B. Hugueney, Y. Lechevallier, F. Rossi, Exploratory analysis of functional data via clustering and optimal segmentation, Neurocomputing, 73(7-9), 1125-1141, (2010)
11. G. M. James, C. A. Sugar, Clustering for Sparsely Sampled Functional Data, Journal of the American Statistical Association, 98(462), 397-408, (2003)
12. X. Liu and M. C. K. Yang, Simultaneous curve registration and clustering for functional data, Computational Statistics and Data Analysis, 53 (4), 1361-1376, (2009)
13. G. McLachlan and D. Peel, Finite Mixture Models, Wiley, (2000)
14. S. K. Ng, G. J. McLachlan, K. Wang, L. Ben-Tovim Jones and S.-W. Ng, A Mixture model with random-effects components for clustering correlated gene-expression profiles, Bioinformatics 22 (14), 1745-1752, (2006)
15. J. O. Ramsay and B. Silverman, Functional data analysis, Springer, New York, (1997)
16. G. Schwarz, Estimating the number of components in a finite mixture model, Annals of Statistics, 6, 461-464, (1978)
17. J. Q. Shi, R. Murray-Smith, D. M. Titterington, Curve Prediction and Clustering with Mixtures of Gaussian Process Functional Regression Models, Statistics and Computing 18 (3), 1573-1375, (2008)
18. C. S. Wong and W. K. Li, On a mixture autoregressive model, Journal of the Royal Statistical Society, Series B, 62, 1, 95-115, (2000)
19. Y. Xiong, D.-Y. Yeung, Time series clustering with ARMA mixtures, Pattern Recognition, 37 (8), 1675-1689, (2004)

# A Convexity of the set $E_{k\ell}$

The set $E_{k\ell}$ defined by:

$$E_{k\ell} = \left\{ t \in [t_1; t_m] \ / \ \pi_{k\ell}(t; \boldsymbol{\alpha}_k) = \max_{1 \leq h \leq L} \pi_{kh}(t; \boldsymbol{\alpha}_k) \right\}.$$

is a convex set of $\mathbb{R}$. In fact, we have the following equalities:

$$
\begin{aligned}
E_{k\ell} &= \left\{ t \in [t_1; t_m] \ / \ \pi_{k\ell}(t; \boldsymbol{\alpha}_k) = \max_{1 \leq h \leq L} \pi_{kh}(t; \boldsymbol{\alpha}_k) \right\} \\
&= \left\{ t \in [t_1; t_m] \ / \ \pi_{kh}(t; \boldsymbol{\alpha}_k) \leq \pi_{k\ell}(t; \boldsymbol{\alpha}_k) \ \text{ for } h = 1, \ldots, L \right\} \\
&= \bigcap_{1 \leq h \leq L} \left\{ t \in [t_1; t_m] \ / \ \pi_{kh}(t; \boldsymbol{\alpha}_k) \leq \pi_{k\ell}(t; \boldsymbol{\alpha}_k) \right\} \\
&= \bigcap_{1 \leq h \leq L} \left\{ t \in [t_1; t_m] \ / \ \ln \frac{\pi_{kh}(t; \boldsymbol{\alpha}_k)}{\pi_{k\ell}(t; \boldsymbol{\alpha}_k)} \leq 0 \right\}
\end{aligned}
$$

From the definition of $\pi_{k\ell}(t; \boldsymbol{\alpha}_k)$ (see equation 6), it can be easily verified that $\ln \frac{\pi_{kh}(t; \boldsymbol{\alpha}_k)}{\pi_{k\ell}(t; \boldsymbol{\alpha}_k)}$ is a linear function of $t$. Consequently, $E_{k\ell}$ is convex, as the intersection of convexes parts of $\mathbb{R}$.