




Clustering multivariate longitudinal observations: The contaminated Gaussian hidden Markov model

Antonio Punzo & Antonello Maruotti

To cite this article: Antonio Punzo & Antonello Maruotti (2015): Clustering multivariate longitudinal observations: The contaminated Gaussian hidden Markov model, Journal of Computational and Graphical Statistics, DOI: [10.1080/10618600.2015.1089776](https://doi.org/10.1080/10618600.2015.1089776)

To link to this article: <http://dx.doi.org/10.1080/10618600.2015.1089776>

 View supplementary material 

 Accepted online: 29 Sep 2015.

 Submit your article to this journal 

 Article views: 2

 View related articles 

 View Crossmark data 

Clustering multivariate longitudinal observations: The contaminated Gaussian hidden Markov model

Antonio Punzo

Department of Economics and Business,

University of Catania

and

Antonello Maruotti

Southampton Statistical Sciences Research Institute,

University of Southampton

Department of Economic, Political Sciences and Modern Languages,

Libera Università Maria Ss. Assunta

Abstract

The Gaussian hidden Markov model (HMM) is widely considered for the analysis of heterogeneous continuous multivariate longitudinal data. To robustify this approach with respect to possible elliptical heavy-tailed departures from normality, due to the presence of outliers, spurious points, or noise (collectively referred to as *bad points* herein), the contaminated Gaussian HMM is here introduced. The contaminated Gaussian distribution represents an elliptical generalization of the Gaussian distribution and allows for automatic detection of bad points in the same natural way as observations are typically assigned to the latent states in the HMM context. Once the model is fitted, each observation has a posterior probability of belonging to a particular state and, inside each state, of being a bad point or not. In addition to the parameters of the classical Gaussian HMM, for each state we have two more parameters, both with a specific and useful interpretation: one controls the proportion of bad points and one specifies their degree of atypicality. A sufficient condition for the identifiability of the model is given, an expectation-conditional maximization algorithm is outlined for parameter estimation and various operational issues are discussed. Using a large scale simulation study, but also an illustrative artificial dataset, we demonstrate the effectiveness of the proposed model in comparison with HMMs of different elliptical distributions, and we also evaluate the performance of some well-known information criteria in selecting the true number of latent states. The model is finally used to fit data on criminal activities in Italian provinces.

Keywords: Robust model-based clustering, Expected-conditional maximization (ECM) algorithm, Model selection, Elliptical distributions, Atypical data

1 Introduction

Hidden Markov models (HMMs) are the state of the art in the analysis of time-dependent data. HMMs have been applied in time series analysis for more than four decades (Baum and Petrie, 1966) and, more recently, in the longitudinal setting (see Bartolucci et al., 2013, Maruotti, 2011, and the reference therein). Serial dependence and heterogeneity in sample units characterize the longitudinal setting and can be properly investigated and accounted for in a HMM framework. Being dependent mixture models, HMMs allow the unambiguously recover of the structure of the data by rigorously defining homogeneous latent subgroups and, simultaneously, provide meaningful interpretation of the inferred partition. For multivariate continuous data, attention is commonly focused on Gaussian HMMs (Holzmann and Schwaiger, 2015; Volant et al., 2014; Bartolucci and Farcomeni, 2010), with few notable exceptions (Lagona et al., 2015; Bulla et al., 2012; Bartolucci and Farcomeni, 2009).

Unfortunately, real data are often contaminated by outliers, spurious points or noise (collectively referred to as *bad points* herein, as in Aitkin and Wilson, 1980) that may affect parameters estimates and the recovering of the latent structure. The attempt of robustly estimating mixture models parameters has led to a heterogeneous literature that includes: noise approaches (Banfield and Raftery, 1993; Fraley and Raftery, 2002), i.e. methods aiming at identifying a noise component (modelled assuming a uniform component-specific distribution), while simultaneously clustering non-noise observations; distance approaches (Rousseeuw and Leroy, 2005; Cerioli, 2010; Garcia-Escudero et al., 2015); distribution-based robust approaches (Peel and McLachlan, 2000; Andrews and McNicholas, 2012). While all these methods offer important contributions to the topic, the last two methods do not allow for the direct detection of bad points. Approaches considering the uniform distribution, if used for discriminant analysis, cannot recognize a new bad observation (an observation that has not been used to fit the model) if it lies outside the support defined by the fitted uniform distribution(s). An alternative to these methods aiming at identifying outlying observations which deviate from the cluster-specific distribution has been recently proposed by Evans et al. (2015). Despite the wide literature on robust estimation of mixture models, there are not many papers dealing with robustness issues in HMMs. In the univariate case, Bulla (2011) introduces a structured HMM to account for outliers in financial time series, Humburg et al. (2008) propose the

use of the t distribution, and Maruotti (2014) considers a bi-square scale estimator in a regression setting. In the multivariate case, Farcomeni and Greco (2015) introduce a robust S-estimator and Bernardi et al. (2014) propose the use of the multivariate t distribution for multivariate financial (time-series) data in a HMM framework.

In this paper, we extend this branch of literature by introducing a joint approach to time-varying robust clustering and bad points detection under a longitudinal setting, extending the standard HMM framework (see Section 2.1). As emphasized by Davies and Gather (1993, see also Hennig, 2002), bad observations should be defined with respect to a reference distribution. Accordingly, the region of bad points can be defined, e.g., as a region where the density of the reference distribution is low. In analogy with other distribution-based approaches (as those based on the t distribution), we choose the Gaussian distribution as the reference distribution but, differently from the t HMM, we replace the multivariate Gaussian state-dependent distribution with a two-component Gaussian mixture (Tukey, 1960) where one (reference) component represents the data we would expect from the given state (i.e. good points) while the other component clusters the bad points; the latter component has a small prior probability, the same component-specific mean and an inflated covariance matrix. Its investigation and use in a clustering framework is still in infancy, although some results have been recently obtained by Punzo and McNicholas (2014a, 2015, 2014b) in a cross-sectional setting. This change makes the model much more robust and allows for automatic detection of bad points. With respect to the latter issue, as it will be better explained later, once the contaminated Gaussian HMM is fitted to the observed longitudinal data, by means of maximum *a posteriori* probabilities, each observation can be first assigned to one of the states and then classified as good or bad; thus, we have a model for simultaneous robust clustering and detection of atypical observations in a longitudinal context. Of course, this is not the only attempt to deal with clustering under a longitudinal setting. Our proposal is somehow related to the models proposed by De la Cruz-Mesia et al. (2008), who introduce a (univariate) hierarchical mixture model, and by McNicholas and Murphy (2010), who consider a (univariate) mixture model in which a decomposed covariance structure is introduced to explicitly account for the relationship between measurements at different time points. However, none of the aforementioned approaches allows for time-varying clustering neither of bad points detection, and, moreover, both have been introduced in a univariate setting

only.

After establishing a sufficient condition for the identifiability of the model (see Section 2.2), in Section 2.3 we outline an *ad hoc* version of the expectation-conditional maximization (ECM) algorithm to estimate model parameters, extending the Baum-Welch iterative procedure (Baum et al., 1970) to deal with contaminated Gaussian distributions. Further operational aspects are discussed in Section 3. In Section 4.1, we illustrate the proposal by a large scale simulation study in order to investigate the empirical behavior of the proposed approach with respect to several factors – such as the number of observed units and times, and the nature of bad points – and in comparison with HMMs of different elliptical distributions. Indeed, different aspects of robustness are going to be described and analyzed. We will consider heavy tails (conditional) distributions as data generation processes as well as distributions with contaminated units in the data. Furthermore, in Section 4.2, we provide insights on information criteria performances in this framework. At last, after an illustration on artificial data (see Section 5.1), we illustrate the proposal in Section 5.2 by analyzing a longitudinal dataset of Italian provinces on which four different crimes rates have been measured from 2005 to 2009, previously analysed in a different context by Viroli (2011). Provinces are clustered and bad points are automatically detected.

2 Methodology

2.1 The model

Let $\{Y_{it}; i = 1, \dots, I, t = 1, \dots, T\}$ denote sequences of multivariate longitudinal observations recorded on I units and T times, where $Y_{it} = (Y_{it1}, \dots, Y_{itP})' \in \mathbb{R}^P$, and let $\{S_{it}; i = 1, \dots, I, t = 1, \dots, T\}$ be a first-order Markov chain defined on the state space $\{1, \dots, k, \dots, K\}$. A HMM is a particular kind of dependent mixture. It is a stochastic process consisting of two parts: the underlying unobserved process $\{S_{it}\}$, fulfilling the Markov property, i.e.

$$\Pr(S_{it} = s_{it} \mid S_{i1} = s_{i1}, S_{i2} = s_{i2}, \dots, S_{it-1} = s_{it-1}) = \Pr(S_{it} = s_{it} \mid S_{it-1} = s_{it-1}),$$

and the state-dependent observation process $\{\mathbf{Y}_{it}\}$ for which the conditional independence property holds, i.e.

$$f(\mathbf{Y}_{it} = \mathbf{y}_{it} \mid \mathbf{Y}_{i1} = \mathbf{y}_{i1}, \dots, \mathbf{Y}_{i,t-1} = \mathbf{y}_{i,t-1}, S_{i1} = s_{i1}, \dots, S_{it} = s_{it}) = f(\mathbf{Y}_{it} = \mathbf{y}_{it} \mid S_{it} = s_{it}),$$

where $f(\cdot)$ is a generic probability density function.

The hidden Markov chain has K states with initial probabilities $\pi_{ik} = \Pr(S_{i1} = k)$, $k = 1, \dots, K$, and transition probabilities

$$\pi_{i,k|j} = \Pr(S_{it} = k \mid S_{i,t-1} = j), \quad t = 2, \dots, T \text{ and } j, k = 1, \dots, K. \quad (1)$$

In (1), k refers to the current state, whereas j refers to the one previously visited; this convention will be used throughout the paper. In the following, for simplicity of explanation, we will consider a homogeneous HMM, that is $\pi_{i,k|j} = \pi_{k|j}$ and $\pi_{ik} = \pi_k$, $i = 1, \dots, I$. Such an assumption can be easily relaxed to include covariates and/or unit-specific random effects as described in Maruotti and Rocci (2012). Thus, we collect the initial probabilities in the K -dimensional vector $\boldsymbol{\pi}$, whereas the time-homogeneous transition probabilities are collected in the $K \times K$ transition matrix $\mathbf{\Pi}$. The conditional density for the observed process is given by a contaminated Gaussian distribution, that is

$$\phi(\mathbf{y}_{it} \mid S_{it} = k; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k, \alpha_k, \eta_k) = \alpha_k \mathcal{N}_P(\mathbf{y}_{it} \mid S_{it} = k; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) + (1 - \alpha_k) \mathcal{N}_P(\mathbf{y}_{it} \mid S_{it} = k; \boldsymbol{\mu}_k, \eta_k \boldsymbol{\Sigma}_k),$$

where $\mathcal{N}_P(\cdot; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$ denotes the P -variate Gaussian distribution with mean $\boldsymbol{\mu}_k$ and covariance matrix $\boldsymbol{\Sigma}_k$, $\alpha_k \in (0, 1)$ is the proportion of good points in state k , and $\eta_k > 1$ is an inflation parameter in state k ; the latter parameter can be also meant as a sort of “degree of atypicality” of the bad point(s), i.e as a measure of how different atypical observations are from the bulk of the (clustered) data.

2.2 Identifiability

An important issue in dealing with the proposed model is to establish its identifiability. Identifiability is a necessary requirement, *inter alia*, for the usual asymptotic theory to hold for maximum likelihood (ML) estimation of the model parameters.

For HMMs, whose state-dependent distributions are assumed to belong to some parametric family, Leroux (1992) shows that identifiability up to label switching follows from identifiability of the marginal finite mixtures (cf. Dannemann et al., 2014, Section 2). In our case, the parametric family is constituted by contaminated Gaussian distributions and the marginal finite mixtures are represented by the finite mixtures of contaminated Gaussian distributions introduced by Punzo and McNicholas (2015). These authors also provide a sufficient condition for the identifiability of their mixture (cf. Punzo and McNicholas, 2015, Proposition 1) that can be summarized as follows:

If $k \neq k_1$ implies

$$\|\boldsymbol{\mu}_k - \boldsymbol{\mu}_{k_1}\|_2^2 + \|\boldsymbol{\Sigma}_k - a\boldsymbol{\Sigma}_{k_1}\|_2^2 \neq 0$$

for all $a > 0$, where $\|\cdot\|_2$ is the Froebenius norm, then a finite mixture of contaminated Gaussian distributions is identifiable.

Accordingly, a finite mixture of contaminated Gaussian distributions is identifiable if two of the K Gaussian distributions representing the good observations have distinct component means and/or non-proportional component covariance matrices. Based on Leroux (1992), the same sufficient condition for identifiability is inherited by our contaminated Gaussian HMM.

2.3 Maximum likelihood estimation

In order to perform ML estimation of the proposed model on the basis of the sample $\{\mathbf{y}_{it}; i = 1, \dots, I, t = 1, \dots, T\}$, the need arises of computing

$$\mathcal{L}(\boldsymbol{\vartheta}) = \prod_{i=1}^I \mathcal{L}_i(\boldsymbol{\vartheta}) = \prod_{i=1}^I \boldsymbol{\pi}' \boldsymbol{\phi}(\mathbf{y}_{i1}) \boldsymbol{\Pi} \boldsymbol{\phi}(\mathbf{y}_{i2}) \boldsymbol{\Pi} \cdots \boldsymbol{\phi}(\mathbf{y}_{iT-1}) \boldsymbol{\Pi} \boldsymbol{\phi}(\mathbf{y}_{iT}) \mathbf{1}_K, \quad (2)$$

where $\boldsymbol{\vartheta} = \{\boldsymbol{\pi}, \boldsymbol{\Pi}, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k, \alpha_k, \eta_k; k = 1, \dots, K\}$ corresponds to the set of all model parameters, $\mathbf{1}_K$ denotes a vector of K ones, and $\boldsymbol{\phi}(\mathbf{y}_{it})$ denotes a $K \times K$ diagonal matrix with conditional densities $\phi(\mathbf{Y}_{it} = \mathbf{y}_{it} | S_{it} = k)$ on the main diagonal. Finding the value of the parameters $\boldsymbol{\vartheta}$ that maximizes the log-transformation of (2) under the constraints $\sum_{k=1}^K \pi_k = 1$, $\sum_{k=1}^K \pi_{k|j} = 1$, $\alpha_k \in (0, 1)$ and $\eta_k > 1$, $k, j = 1, \dots, K$, is not an easy problem since (2) is not available in an analytically convenient form. Efficient computation of (2) may be performed by exploiting a forward recursion described in the HMM literature (see, e.g., Zucchini and MacDonald, 2009).

In this relatively general framework, an expectation-conditional maximization (ECM) algorithm (Meng and Rubin, 1993) is used for fitting our contaminated model. The ECM is a variant of the classical EM algorithm (Baum et al., 1970; Dempster et al., 1977), which is a natural approach for ML estimation when data are incomplete. In our setting, there are two sources of missing data: one arises from the fact that we do not know state membership and its evolution over time, and the other from the fact that we do not know whether an observation clustered in a specific state is good or bad.

Formally, let us define the unobserved state membership $\mathbf{z}_{it} = (z_{it1}, \dots, z_{itk}, \dots, z_{itK})'$, the unobserved states transition

$$\mathbf{zz}_{it} = \begin{pmatrix} zz_{it11} & \cdots & zz_{it1k} & \cdots & zz_{it1K} \\ \vdots & & \vdots & & \vdots \\ zz_{itj1} & \cdots & zz_{itjk} & \cdots & zz_{itjK} \\ \vdots & & \vdots & & \vdots \\ zz_{itK1} & \cdots & zz_{itKk} & \cdots & zz_{itKK} \end{pmatrix},$$

and the unobserved state-specific membership to the good points $\mathbf{v}_{it} = (v_{it1}, \dots, v_{itk}, \dots, v_{itK})'$, as missing data, with

$$z_{itk} = \begin{cases} 1 & \text{if } S_t = k \\ 0 & \text{otherwise} \end{cases}, \quad zz_{itjk} = \begin{cases} 1 & \text{if } S_{it-1} = j \text{ and } S_{it} = k \\ 0 & \text{otherwise} \end{cases},$$

and $v_{itk} = 1$ if observation i at time t in state k is a good point and $v_{itk} = 0$ if it is a bad point. Therefore, the complete data are given by $\mathcal{C} = \{\mathbf{y}_{it}, \mathbf{z}_{it}, \mathbf{zz}_{it}, \mathbf{v}_{it}; i = 1, \dots, I, t = 1, \dots, T\}$ and the complete-data log-likelihood can be written as

$$\ell_c(\boldsymbol{\theta} | \mathcal{C}) = \ell_{c_1}(\boldsymbol{\pi} | \mathcal{C}) + \ell_{c_2}(\mathbf{\Pi} | \mathcal{C}) + \ell_{c_3}(\boldsymbol{\alpha} | \mathcal{C}) + \ell_{c_4}(\boldsymbol{\mu}, \boldsymbol{\Sigma}, \boldsymbol{\eta} | \mathcal{C})$$

where

$$\begin{aligned}
 \ell_{c_1}(\boldsymbol{\pi} | \mathcal{C}) &= \sum_{i=1}^I \sum_{k=1}^K z_{i1k} \log(\pi_k) \\
 \ell_{c_2}(\boldsymbol{\Pi} | \mathcal{C}) &= \sum_{i=1}^I \sum_{t=2}^T \sum_{k=1}^K \sum_{j=1}^K z_{itjk} \log(\pi_{k|j}) \\
 \ell_{c_3}(\boldsymbol{\alpha} | \mathcal{C}) &= \sum_{i=1}^I \sum_{t=1}^T \sum_{k=1}^K z_{itk} [v_{itk} \log(\alpha_k) + (1 - v_{itk}) \log(1 - \alpha_k)] \\
 \ell_{c_4}(\boldsymbol{\mu}, \boldsymbol{\Sigma}, \boldsymbol{\eta} | \mathcal{C}) &= -\frac{1}{2} \sum_{i=1}^I \sum_{t=1}^T \sum_{k=1}^K \left[z_{itk} \log |\boldsymbol{\Sigma}_k| + P z_{itk} (1 - v_{itk}) \log(\eta_k) \right. \\
 &\quad \left. + z_{itk} \left(v_{itk} - \frac{1 - v_{itk}}{\eta_k} \right) (\mathbf{y}_{it} - \boldsymbol{\mu}_k)' \boldsymbol{\Sigma}_k^{-1} (\mathbf{y}_{it} - \boldsymbol{\mu}_k) \right],
 \end{aligned}$$

where $\boldsymbol{\mu} = \{\boldsymbol{\mu}_k; k = 1, \dots, K\}$, $\boldsymbol{\Sigma} = \{\boldsymbol{\Sigma}_k; k = 1, \dots, K\}$, $\boldsymbol{\eta} = \{\eta_k; k = 1, \dots, K\}$, and $\boldsymbol{\alpha} = \{\alpha_k; k = 1, \dots, K\}$.

The E-step, at the $(r + 1)$ -th iteration, computes the conditional expectations of ℓ_c with respect to $\boldsymbol{\theta}$, given the observed data and the current estimates of the parameters. To do this, we replace z_{itk} and z_{itjk} with their conditional expectations, namely, $\tilde{z}_{itk}^{(r)}$ and $\tilde{z}_{itjk}^{(r)}$ (for computational details, see Section ?? in the Supplementary Material) and v_{itk} with

$$\tilde{v}_{itk}^{(r)} = E(V_{itk} | \mathbf{y}_{it}, \boldsymbol{\theta}^{(r)}) = \frac{\alpha_k \mathcal{N}_P(\mathbf{y}_{it} | S_{it} = k; \boldsymbol{\mu}_k^{(r)}, \boldsymbol{\Sigma}_k^{(r)})}{\phi(\mathbf{y}_{it} | S_{it} = k; \boldsymbol{\mu}_k^{(r)}, \boldsymbol{\Sigma}_k^{(r)}, \alpha_k^{(r)}, \eta_k^{(r)})}, \quad (3)$$

where V_{itk} is the random variable related to v_{itk} .

At the first CM-step of the $(r + 1)$ -th iteration, maximizing with respect to π_k , $\mathbf{\Pi}$, $\boldsymbol{\mu}_k$, $\boldsymbol{\Sigma}_k$ and α_k yields

$$\begin{aligned} \pi_k^{(r+1)} &= \frac{\sum_{i=1}^I \tilde{z}_{i1k}^{(r)}}{I}, & \pi_{klj}^{(r+1)} &= \frac{\sum_{i=1}^I \sum_{t=2}^T \tilde{z}_{itjk}^{(r)}}{\sum_{i=1}^I \sum_{t=2}^T \sum_{k=1}^K \tilde{z}_{itjk}^{(r)}}, & \alpha_k^{(r+1)} &= \frac{\sum_{i=1}^I \sum_{t=1}^T \tilde{z}_{itk}^{(r)} \tilde{v}_{itk}^{(r)}}{\sum_{i=1}^I \sum_{t=1}^T \tilde{z}_{itk}^{(r)}}, \\ \boldsymbol{\mu}_k^{(r+1)} &= \frac{\sum_{i=1}^I \sum_{t=1}^T \tilde{z}_{itk}^{(r)} \left(\tilde{v}_{itk}^{(r)} + \frac{1 - \tilde{v}_{itk}^{(r)}}{\eta_k^{(r)}} \right) \mathbf{y}_{it}}{\sum_{i=1}^I \sum_{t=1}^T \tilde{z}_{itk}^{(r)} \left(\tilde{v}_{itk}^{(r)} + \frac{1 - \tilde{v}_{itk}^{(r)}}{\eta_k^{(r)}} \right)}, \end{aligned} \quad (4)$$

$$\boldsymbol{\Sigma}_k^{(r+1)} = \frac{\sum_{i=1}^I \sum_{t=1}^T \tilde{z}_{itk}^{(r)} \left(\tilde{v}_{itk}^{(r)} + \frac{1 - \tilde{v}_{itk}^{(r)}}{\eta_k^{(r)}} \right) (\mathbf{y}_{it} - \boldsymbol{\mu}_k) (\mathbf{y}_{it} - \boldsymbol{\mu}_k)'}{\sum_{i=1}^I \sum_{t=1}^T \tilde{z}_{itk}^{(r)}}. \quad (5)$$

At the second CM-step of the $(r + 1)$ -th iteration, we maximize the expectation of the complete-data log-likelihood with respect to η_k , fixing all other parameters to their estimated values at the first CM-step. In particular, we have to maximize

$$-\frac{P}{2} \sum_{i=1}^I \sum_{t=1}^T \tilde{z}_{itk}^{(r)} (1 - \tilde{v}_{itk}^{(r)}) \log(\eta_k) - \frac{1}{2} \sum_{i=1}^I \sum_{t=1}^T \tilde{z}_{itk}^{(r)} \frac{1 - \tilde{v}_{itk}^{(r)}}{\eta_k} (\mathbf{y}_{it} - \boldsymbol{\mu}_k^{(r+1)})' (\boldsymbol{\Sigma}_k^{(r+1)})^{-1} (\mathbf{y}_{it} - \boldsymbol{\mu}_k^{(r+1)})$$

with respect to η_k , under the constraint $\eta_k > 1$, $k = 1, \dots, K$. As a closed form solution is not analytically available, the `optimize()` function in the `stats` package of the R software (R Core Team, 2013) is used to perform numerical search of the maximum of the previous expression.

3 Operational aspects

3.1 Note on robustness

Let us focus on the weights $\left(\tilde{v}_{itk} + \frac{1-\tilde{v}_{itk}}{\eta_k}\right)$ in (4) and (5). We can rewrite \tilde{v}_{itk} as an explicit function of the squared Mahalanobis distance, say δ , as

$$h(\delta; \alpha_k, \eta_k) = \frac{\alpha_k \exp\left(-\frac{\delta}{2}\right)}{\alpha_k \exp\left(-\frac{\delta}{2}\right) + \frac{(1-\alpha_k)}{\sqrt{\eta_k}} \exp\left(-\frac{\delta}{2\eta_k}\right)} = \frac{1}{1 + \frac{(1-\alpha_k)}{\alpha_k} \frac{1}{\sqrt{\eta_k}} \exp\left[\frac{\delta}{2}\left(1 - \frac{1}{\eta_k}\right)\right]},$$

with $\delta \geq 0$.

Due to $\eta_k > 1$, $h(\delta; \alpha_k, \eta_k)$ is a decreasing function of δ . Accordingly,

$$w(\delta; \alpha_k, \eta_k) = \left(\tilde{v}_{itk} + \frac{1-\tilde{v}_{itk}}{\eta_k}\right) = h(\delta; \alpha_k, \eta_k) + \frac{1-h(\delta; \alpha_k, \eta_k)}{\eta_k} = \frac{1}{\eta_k} [1 + (\eta_k - 1)h(\delta; \alpha_k, \eta_k)],$$

which is an increasing function of $h(\delta; \alpha_k, \eta_k)$ and, thus, a decreasing function of δ . Therefore $w(\delta; \alpha_k, \eta_k)$ reduces the effect of bad points in the estimation of $\boldsymbol{\mu}_k$ and $\boldsymbol{\Sigma}_k$, $k = 1, \dots, K$, so providing a robust way to estimate these parameters (see also Punzo and McNicholas, 2015).

3.2 Detection of bad points and further constraints

For the proposed model, the classification of an observation \mathbf{y}_{it} is a two-step procedure

Step 1. determine state membership via local or global decoding procedures (see Section ?? in the Supplementary Material);

Step 2. establish if it is either a good or a bad observation in that state.

Once the hidden path is inferred, for each observation, we look at \tilde{v}_{itk} for the inferred state and \mathbf{y}_{it} is good if $\tilde{v}_{itk} > 0.5$ and bad otherwise.

Bearing in mind that $(1 - \alpha_k)$ represents the proportion of bad points and η_k denotes the degree of contamination, we would require that in the k th hidden state, $k = 1, \dots, K$, the proportion of good data is at least equal to a fixed value α_k^* . In this case, the `optimize()` function is also used

for a numerical search of the maximum $\alpha_k^{(r+1)}$, over the interval $(\alpha_k^*, 1)$, of the function

$$\sum_{i=1}^I \sum_{t=1}^T \tilde{z}_{itk}^{(r)} \left[\tilde{v}_{itk}^{(r)} \log \alpha_k + (1 - \tilde{v}_{itk}^{(r)}) \log (1 - \alpha_k) \right].$$

In both the simulation study and the empirical application, we use this approach to update α_k and we take $\alpha_k^* = 0.5$, for $k = 1, \dots, K$. The value 0.5 is justified because, in robust statistics, it is usually assumed that at least half of the points are good (cf. Hennig, 2002, p. 250). Note that it is also possible to fix α_k and/or η_k *a priori*. This is somewhat analogous to the trimmed clustering approach, where one must to specify the proportion of outliers (the so-called trimming proportion) in advance (cf. Fritz et al., 2012).

4 Simulation studies

In this section we investigate various aspects of the proposed model through large-scale simulation studies performed using R (R Core Team, 2013).

4.1 Comparison between HMMs of elliptical distributions

The first simulation study aims to demonstrate the effectiveness of the proposed model in comparison with HMMs of some elliptical distributions. A general feedback on advantages and drawbacks of each model is also given. We compare: the HMM of Gaussian distributions (NHMM); the HMM of t distributions (t HMM); the HMM of contaminated Gaussian distributions (CNHMM). To generate the data, we consider the following five data generation processes with bivariate ($P = 2$) state-specific distributions and $K = 2$ hidden states:

- a) NHMM;
- b) t HMM with $\nu_1 = 4$ and $\nu_2 = 10$ degrees of freedom;
- c) CNHMM with $\alpha_1 = 0.9$, $\alpha_2 = 0.8$, $\eta_1 = 2$, and $\eta_2 = 20$;
- d) NHMM with 1% of points randomly substituted by high atypical points with coordinates $(0, y_{it2}^*)$, where y_{it2}^* is generated from a uniform distribution over the interval $(10, 15)$.

- e) NHMM with 5% of points randomly substituted by noise points generated from a uniform distribution over the interval $(-10, 10)$ on each dimension.

All of these data generation processes share the following common parameters

$$\pi_1 = 0.3, \mathbf{\Pi} = \begin{pmatrix} 0.8 & 0.2 \\ 0.2 & 0.8 \end{pmatrix}, \boldsymbol{\mu}_1 = \begin{pmatrix} 0 \\ -3 \end{pmatrix}, \boldsymbol{\mu}_2 = -\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1 = \begin{pmatrix} 1 & -0.5 \\ -0.5 & 1 \end{pmatrix}, \text{ and } \boldsymbol{\Sigma}_2 = \begin{pmatrix} 1 & 0.5 \\ 0.5 & 1 \end{pmatrix}.$$

The five scenarios above cover different situations which may arise dealing with real-world data: no bad points for scenario *a*), heavy-tails conditional distributions for scenarios *b*) and *c*), and two different types of bad points for scenarios *d*) and *e*). Under each scenario, we simulate 100 samples considering two experimental factors: the number of analyzed units I (50, 100, and 200) and the number of repeated measurements T (5, 10, and 20). This yields a total of 4,500 generated datasets. On each generated dataset, the EM-based algorithm of the three competing models is directly run with $K = 2$, is initialized according to the partition provided by the K -means method, and is stopped when the difference between the updated parameter estimates of two consecutive iterations is less than 10^{-4} .

For comparison's sake, we report the bias (BIAS) and the standard deviation (STD) of the estimates for the initial weight π_1 , the transition probabilities $\pi_{1|1}$ and $\pi_{2|2}$ (diagonal elements of $\mathbf{\Pi}$), the univariate means μ_{11} and μ_{21} (elements of $\boldsymbol{\mu}_1$), and the univariate means μ_{12} and μ_{22} (elements of $\boldsymbol{\mu}_2$). We would remark that HMMs, and mixture models in general, are affected by label switching issues (see, e.g., Yao, 2012), which render estimators evaluation using simulations more complex. There are no generally accepted labeling methods. In our simulation study, because of true values $\pi_1 = 0.3$ and $\pi_2 = 0.7$, we simply attribute the label 1 to the state with the lowest estimated initial probability.

The obtained results are reported in Tables ??-?? in the Supplementary Material. As concerns the estimates of $\boldsymbol{\mu}_1$ and $\boldsymbol{\mu}_2$, we note the following general findings. Under scenario *a*), that is when there are no bad points, all the approaches perform comparably, as expected since, in this situation, the t HMM and the CNHMM tend to the NHMM. Under scenarios *b*) and *c*), the robust approaches, t HMM and CNHMM, are better than the traditional NHMM, especially when data are generated by the CNHMM (see, e.g., Table ?? in the Supplementary Material). Moreover, the fitted t HMM

and CNHMM perform comparably in both scenarios; such a comparable behavior agrees with the simulation results of Little (1988) about the t and the contaminated Gaussian distributions. Even under scenarios d) and e), the t HMM and the CNHMM perform comparably and much better than the NHMM. In particular, under scenario d), we note how the NHMM-estimates of the means μ_{21} and μ_{22} for the second dimension are mainly affected by the bad points; this is a natural result if we recall that these points are bad due to the value of the second dimension.

Referring to hidden parameters, as expected, under scenario a) all the considered approaches perform well, providing unbiased estimates for both the transition probability matrix $\mathbf{\Pi}$ and the initial probabilities π_1 and π_2 . Even under scenario b), the three approaches performs well and almost comparably, with a slight worse performance for the NHMM. Under scenario c), the traditional NHMM underestimates the first element $\pi_{1|1}$ on the main diagonal of the transition probability matrix of the hidden chain. In other words, this approach estimates a higher number of transitions, from state 1 to state 2, than the ones assumed by the model used to generate the data. Moreover, the initial probability π_1 of state 1 is slightly underestimated too. Of course, these could represent issues if the underlying latent structure is of interest. On the contrary, under the same scenario, the t HMM and the CNHMM perform well, providing unbiased estimates for $\mathbf{\Pi}$, π_1 , and π_2 . The findings about scenario d) are similar to those under scenario b). Finally, scenario e), that is the noise case, is the most problematic for the traditional NHMM. In fact, being the noisy observations drawn at random from a Bernoulli distribution with parameter equal to 0.05, they may produce sudden changes in the latent structure, altering its dynamics. In particular, the initial probability π_1 , as well as the the first element $\pi_{1|1}$ on the main diagonal of the transition probability matrix of the hidden chain, are strongly underestimated. An *a posteriori* analysis, as well as the magnitude of the bias for π_1 , reveals that the NHMM identifies a state (i.e. the persistent state) where all the good observations are clustered and another state (i.e. the non-persistent state) where all the noisy observations are grouped. The resulting hidden structure is, thus, characterized by sudden changes towards the non-persistent state followed by successive changes to the persistent state.

Table ?? in the Supplementary Material summarizes the obtained average misclassification rates. Misclassification rates are computed via the `classError()` function of the **mclust** package for R (Fraley et al., 2015). Note that, under scenarios d) and e), misclassification rates are computed

only with respect to the true good observations. The results in Table ?? corroborate the previous simulation findings; in particular, the robust approaches have a similar behavior and they, apart from scenario *a*), are better than the traditional NHMM, especially under scenarios *c*) and *e*).

Thus far, the *t*HMM and the CNHMM have shown a similar behavior; however, as previously emphasized, the CNHMM has the advantage to allow for the automatic detection of bad points. For the purpose of evaluation of this aspect, we report the true positive rate (TPR), measuring the proportion of bad points that are correctly identified as bad points, and the false positive rate (FPR), corresponding to the proportion of good points incorrectly classified as bad points.

Table 1 reports these measures for scenarios *d*) and *e*).

We note almost optimal results under scenario *d*) and for the FPRs under scenario *e*). Furthermore, the fact that the TPRs do not approach at one under scenario *e*) is not necessarily an error: the way the noisy points are inserted into the data makes possible that some of them will have values related to good points and, as such, these points will be detected as good points by our model.

Finally, to have an idea of the computational burden required by our ECM algorithm, Table 2 shows the average elapsed time (in seconds over 100 replications) to fit a single CNHMM under scenario *c*). Computation is performed on a Windows 8.1 PC, with Intel i7 3.50GHz CPU, 16.0 GB RAM, using R 32 bit, and the elapsed time is computed via the `proc.time()` function of the **base** package. To make the analysis finer, simulations with $I = 150$ and $T = 15$ have been added.

The minimum average elapsed time of 0.918 seconds is obtained in correspondence of the pair $(I = 50, T = 5)$, while the maximum (82.731 seconds) is obtained for the pair $(I = 200, T = 20)$. Furthermore, the average elapsed time seems to be a function of the overall size IT of the data, although the time-length of the panel affects the elapsed time slightly more than the sample size; this conjecture is corroborated by the plot in Figure 1.

4.2 Selecting the number of hidden states

The performance of the information criteria illustrated in Section ?? in the Supplementary Material is here investigated for the CNHMM. To generate the data, we consider the following two CNHMMs with $P = 2$ dimensions:

f) the same two-state model considered under scenario c);

g) a three-state model with the following parameters

$$\pi_1 = 0.16, \quad \pi_2 = 0.34, \quad \mathbf{\Pi} = \begin{pmatrix} 0.8 & 0.1 & 0.1 \\ 0.1 & 0.8 & 0.1 \\ 0.1 & 0.1 & 0.8 \end{pmatrix}, \quad \boldsymbol{\mu}_1 = \begin{pmatrix} 0 \\ -8 \end{pmatrix}, \quad \boldsymbol{\mu}_2 = \begin{pmatrix} 0 \\ 0 \end{pmatrix},$$

$$\boldsymbol{\mu}_3 = \begin{pmatrix} 0 \\ 8 \end{pmatrix}, \quad \boldsymbol{\Sigma}_1 = \begin{pmatrix} 1 & -0.5 \\ -0.5 & 1 \end{pmatrix}, \quad \boldsymbol{\Sigma}_2 = \begin{pmatrix} 0.5 & 0 \\ 0 & 0.5 \end{pmatrix}, \quad \boldsymbol{\Sigma}_3 = \begin{pmatrix} 1 & 0.5 \\ 0.5 & 1 \end{pmatrix},$$

$$\alpha_1 = 0.9, \quad \alpha_2 = 0.8, \quad \alpha_3 = 0.9, \quad \eta_1 = 2, \quad \eta_2 = 5, \quad \text{and} \quad \eta_3 = 10.$$

Under each scenario, we simulate 500 samples considering two experimental factors: the number of analyzed units I (50 and 200) and the number of repeated measurements T (5 and 10). This yields a total of 4,000 generated datasets. On each generated dataset, the ECM algorithm for the CNHMM is run for $K \in \{1, 2, 3, 4, 5\}$, is initialized according to the partition provided by the K -means method, and is stopped when the difference between the updated parameter estimates of two consecutive iterations is less than 10^{-4} .

Table ?? in the Supplementary Material summarizes the obtained results in terms of selection rate (over the 500 replications); the selection rate is defined here as the proportion of times each value of K is selected by the corresponding criterion shown on the top of the column. The rows related to the true value of K are highlighted in gray; to facilitate performance evaluation, the last row of Table ?? gives the mean selection rate of each criterion, computed over the true values of K (i.e., computed over the gray rows).

The BIC and the ICL perform comparably and much better than the AIC that, especially under scenario g), tends to overestimate the number of states.

5 Illustrative examples

5.1 Artificial longitudinal blue crabs data

This section is based on an artificial longitudinal version of the very popular crab dataset of Campbell and Mahon (1974). Attention is focused on the sample of $I = 100$ blue crabs of the genus *Leptograpsus*, subdivided in two groups of equal size ($\pi_1 = \pi_2 = 0.5$). For each specimen, we consider $P = 2$ measurements (in millimeters), namely the rear width (RW) and the length along the midline of the carapace (CL). Mardia's test suggests that it is reasonable to assume that the two group-conditional distributions are bivariate normal (see Greselin et al., 2011, Greselin and Punzo, 2013, and Bagnato et al., 2014 for details). The ML estimates of the parameters $\boldsymbol{\mu}_1$, $\boldsymbol{\mu}_2$, $\boldsymbol{\Sigma}_1$, and $\boldsymbol{\Sigma}_2$ are given in Greselin et al. (2011, p. 158); based on these estimates, and further introducing a transition probabilities matrix

$$\boldsymbol{\Pi} = \begin{pmatrix} 0.9 & 0.1 \\ 0.1 & 0.9 \end{pmatrix},$$

we randomly generate a longitudinal version of this dataset on $T = 5$ times, based on the NHMM; the dataset is available at <http://www.economia.unict.it/punzo/Data.htm>. The scatterplots of the generated data, for each $t \in \{1, \dots, 5\}$, are displayed in Figure 2.

In the fashion of Peel and McLachlan (2000), eight ‘‘perturbed’’ datasets are created by substituting the original value of CL for the 17th point at time 1 (highlighted by a bullet in Figure 2(a)) with eight atypical values shown in the first column of Table 3. We recall that, in the cross-section setting, the aim of Peel and McLachlan (2000) was to show that, unlike Gaussian mixtures, mixtures of t -distributions are robust to these perturbations when applied for clustering.

Ceteris paribus with Peel and McLachlan (2000), we directly fit the NHMM, the t HMM, and the CNHMM, with $K = 2$. For each of the three competing models, Table 3 reports the number of misallocated observations for each perturbed dataset.

It can be seen that, as expected, the t HMM and the CNHMM clusterings are more robust to these perturbations than the NHMM clustering. However, the CNHMM is systematically the most robust to these perturbations, with the number of misallocated observations remaining fixed at 3 regardless of the particular value perturbed; this is in contrast to the NHMM where the number of

misclassifications changes as the extent of the perturbation increases. Interestingly, the CNHMM always detects the perturbed value as a bad point regardless from its magnitude. Furthermore, by recalling that the original value of CL for the 17th point at time $t = 1$ was 32.158, it is also interesting to note that the estimated value of η_k (in the group containing the outlier, which is always the first group) increases as the value of this point further departs from its true value (refer to the fifth column of Table 3). In connection with this aspect, we also report the estimated posterior probability to be a good point (see equation (3)) for the 17th observation at time 1 (refer to the sixth column of Table 3); as we can see, farther the perturbed value is from its group, lower is its probability to be a good point. Such a low probability is also related to the down-weighting of this bad point in the estimation of μ_1 and Σ_1 , and this is an important aspect for the robust estimation of these parameters (see Section 3.1).

Finally note that, in all the considered cases, the CNHMM detects a false positive bad point, which can be easily seen at the bottom-left corner of Figure 2(c).

5.2 Criminal activities in Italian provinces

In this section we analyze data on criminal activities in Italy. Data are taken from an Italian financial newspaper (*Il Sole 24 Ore*, www.ilsole24ore.com), and have been previously analyzed by Viroli (2011). Italian crime has specific features. Firstly, criminal patterns may vary across times and types of activity; secondly, *organized crime* has often territorial roots in specific Italian areas. Bearing this in mind, we would capture differences in (non-violent) criminal activities across time, types of crimes and territorial units, aiming at identifying different levels of safety conditions (represented by the hidden states).

Our analysis focuses on 103 NUTS3 (European Nomenclature of Territorial Units for Statistics) units in Italy, on which we recorded $P = 4$ criminal indicators: home-invasion robberies (per 100,000 residents; HOME); teenage crime rate (per 1,000 residents; TEEN); reported robberies (per 100,000 residents; ROB); rate of muggings and pickpockets (per 100,000 residents; PICK) over five years, from 2005 to 2009. Summaries of the evolution over time of these variables are reported in Figure ?? in the Supplementary Material. We observe an increase in home-invasion robberies over time, whilst all other indicators do not show, at first glance, any significant temporal

variations. Moreover, it is clear, even from simple boxplots, that some units show *unusual values*.

In analyzing the dataset, the most interesting scientific question concerns the existence of areas with similar criminal rates in Italy. Also of interest is the strength of time dependence as measured by the transition probability matrix. Indeed, a strong time-dependence implies no improvements in hindering criminal activities. As described in previous sections, we jointly allow for time-varying clustering as well as for atypical data detections, that may affect the resulting clustering if not properly accounted for. The central idea is that the hidden states cope with the temporal and the spatial structure of the data, and that the contaminated Gaussian distributions can account for atypical data.

On these data, we fit the proposed model with a number of hidden states ranging from 1 to 10. For completeness, we also fit the NHMM and the t HMM. The results are reported in Table 4 in terms of AIC, BIC and ICL. For each value of K , we adopt a K -means approach, with 20 random starting points, to initialize the ECM algorithm, and we report the results corresponding to the best solution in terms of likelihood. On the basis of these results, we conclude that $K = 7$ is a suitable number of hidden states for the considered dataset. This value of K corresponds to the maximum value of both the BIC and ICL criteria, whereas the AIC selects 9 states. However, in the simulation study of Section 4.2, we show that BIC and ICL perform well in recovering the *true* number of hidden states, whereas the AIC may overestimate this number.

On the basis of the estimates of the parameters under the selected CNHMM with $K = 7$ (see Table 5), we conclude that Italian provinces show territorial-specific characteristics and heterogeneous criminal-related situations. As it is clear from the estimated state-specific mean vectors (Table 5) and the inferred clustering structure (Figure ?? in the Supplementary Material), States 1, 4 and 7 are characterized by similar home-invasion robberies rates. They differ in the other indicators. State 1 identifies high teen-crime and low robberies rates, whereas in State 7 home-invasions arise along with reported robberies and pickpockets. State 4, instead, is characterized by home robberies only. With few exceptions, these three states are observed in the most industrialized areas, e.g. in North-North-West provinces. State 5 characterizes big cities (e.g. Rome, Turin and Milan) and touristic places (e.g. Rimini). These are the provinces with the highest values of home robberies, teenage crimes and reported muggings, and therefore the most dangerous ones

in terms of the crime indicators considered in this analysis. Safer provinces are clustered in State 3, whereas unsafe southern provinces, which are notoriously and particularly unsafe in terms of robberies and muggings are clustered in State 6. At last, North-East and Center-North provinces with high rates of pickpockets are clustered in State 2. As we can note by the estimated transition probability matrix

$$\mathbf{\Pi} = \begin{pmatrix} 0.949 & 0.021 & 0.030 & 0.000 & 0.000 & 0.000 & 0.000 \\ 0.000 & 0.984 & 0.000 & 0.000 & 0.000 & 0.000 & 0.016 \\ 0.000 & 0.014 & 0.980 & 0.000 & 0.000 & 0.006 & 0.000 \\ 0.000 & 0.000 & 0.000 & 1.000 & 0.000 & 0.000 & 0.000 \\ 0.000 & 0.000 & 0.000 & 0.000 & 1.000 & 0.000 & 0.000 \\ 0.000 & 0.000 & 0.000 & 0.000 & 0.000 & 1.000 & 0.000 \\ 0.000 & 0.000 & 0.000 & 0.000 & 0.000 & 0.000 & 1.000 \end{pmatrix},$$

transitions between states are rare, persistence is the norm (as expected) and, as a consequence, the global and local decoding procedures provide the same inferred clustering structure. As also discussed in Bulla (2011), the strong persistence is a consequence of the robustification of the HMM. Indeed the NHMM with $K = 7$ has less persistent states. Furthermore, the few transitions are a clear indication that the political action to reduce these criminal activities have not achieved any improving results across the provinces. On the other hand, this very high persistence may indicate the absence of a time-varying clustering structure. The time-varying clustering structure is thus investigated. The data are fitted assuming an additional constraint on the selected model: the identity transition probability matrix is imposed to check for time-constant clustering. The AIC, BIC and ICL values obtained under this model are -18486.29, -18797.19 and -18799.64, respectively; all of them are lower than the values reported in Table 4. Therefore, the time-dependence for the clustering structure is supported by the data.

Other considerations arising from Table 5 are that atypical data can be detected under specific states only. Indeed, the probability of having good data in States 1, 3 and 4 is substantially 1, i.e. all data are estimated as good points. Bad points identified by the CNHMM are displayed in Figure ?? in the Supplementary Material. Naples and Caserta are estimated as atypical in State 6, the one clustering most of Southern provinces, over all the time periods. In Viroli (2011), a state is

devoted to cluster these two provinces only (a similar result is obtained from the t HMM, in which a further state is devoted to cluster Naples and Caserta only, i.e. $K = 8$, and a fuzzier clustering is estimated). Few other provinces are also identified as possible bad points in other states. The η parameters provide the degree of atypicality of these bad points (cf. Section 5.1).

6 Discussion

We have presented a new model for clustering multivariate longitudinal data in a hidden Markov framework, which is robust to outlying observations, spurious observations, or noise, which we collectively referred to as *bad points*. On real and simulated data, we have demonstrated that our model works better than the standard hidden Markov model based on the Gaussian distribution and comparably well with respect to the hidden Markov model based on the heavy-tails t distribution. The main advantage of our model lies in automatic *bad points* detection that is performed by using a maximum *a posteriori* rule. In addition to these advantages, the choice of this approach is motivated by considerable conceptual and computational simplicity in the attempt to generalize the classical Gaussian HMM in terms of robustness and automatic detection of bad points; indeed, only minor modifications to the standard EM algorithm for the Gaussian HMM are involved (cf. Section 2.3). We also investigated information criteria behavior in this framework and observed good BIC and ICL performances in recovering the hidden structure.

There are different possibilities for further work, three of which are worth mentioning. First of all, our approach should be extended to large dimensions. This can be done for instance constraining the covariance matrices of the states, at the price of more complex estimation strategies. As often in the longitudinal setting, covariates information are also available along with multiple response variables. A straightforward extension would deal with the regression framework, in which contamination may arise in the covariate part of the regression model. Similarly, the homogeneous assumption on the hidden Markov chain can be easily relaxed, allowing for time and/or individual specific Markov chains, as well as it is possible to allow for partially missing observations without too much effort, in a missing at random setting.

References

- Aitkin, M. and G. T. Wilson (1980). Mixture models, outliers, and the EM algorithm. *Technometrics* 22(3), 325–331.
- Andrews, J. L. and P. D. McNicholas (2012). Model-based clustering, classification, and discriminant analysis via mixtures of multivariate t -distributions. *Statistics and Computing* 22(5), 1021–1029.
- Bagnato, L., F. Greselin, and A. Punzo (2014). On the spectral decomposition in normal discriminant analysis. *Communications in Statistics - Simulation and Computation* 43(6), 1471–1489.
- Banfield, J. D. and A. E. Raftery (1993). Model-based Gaussian and non-Gaussian clustering. *Biometrics* 49(3), 803–821.
- Bartolucci, F. and A. Farcomeni (2009). A multivariate extension of the dynamic logit model for longitudinal data based on a latent Markov heterogeneity structure. *Journal of the American Statistical Association* 104(486), 816–831.
- Bartolucci, F. and A. Farcomeni (2010). A note on the mixture transition distribution and hidden Markov models. *Journal of Time Series Analysis* 31(2), 132–138.
- Bartolucci, F., A. Farcomeni, and F. Pennoni (2013). *Latent Markov models for longitudinal data*. London: Taylor & Francis.
- Baum, L. E. and T. Petrie (1966). Statistical inference for probabilistic functions of finite state Markov chains. *The Annals of Mathematical Statistics* 37(6), 1554–1563.
- Baum, L. E., T. Petrie, G. Soules, and N. Weiss (1970). A maximization technique occurring in the statistical analysis of probabilistic functions of Markov chains. *The Annals of Mathematical Statistics* 41(1), 164–171.
- Bernardi, M., A. Maruotti, and L. Petrella (2014). Multivariate Markov-switching models and tail risk interdependence. arXiv.org e-print 1312.6407, available at: <http://arxiv.org/abs/1312.6407>.
- Bulla, J. (2011). Hidden Markov models with t components. Increased persistence and other

- aspects. *Quantitative Finance* 11(3), 459–475.
- Bulla, J., F. Lagona, A. Maruotti, and M. Picone (2012). A multivariate hidden Markov model for the identification of sea regimes from incomplete skewed and circular time series. *Journal of Agricultural, Biological and Environmental Statistics* 17(4), 544–567.
- Campbell, N. A. and R. J. Mahon (1974). A multivariate study of variation in two species of rock crab of genus *Leptograpsus*. *Australian Journal of Zoology* 22(3), 417–425.
- Ceroli, A. (2010). Multivariate outlier detection with high-breakdown estimators. *Journal of the American Statistical Association* 105(489), 147–156.
- Dannemann, J., H. Holzmann, and A. Leister (2014). Semiparametric hidden Markov models: identifiability and estimation. *Wiley Interdisciplinary Reviews: Computational Statistics* 6(6), 418–425.
- Davies, L. and U. Gather (1993). The identification of multiple outliers. *Journal of the American Statistical Association* 88(423), 782–792.
- De la Cruz-Mesia, R., A. Quintana, Fernando, and G. Marshall (2008). Model-based clustering for longitudinal data. *Computational Statistics & Data Analysis* 52(3), 1441–1457.
- Dempster, A. P., N. M. Laird, and D. B. Rubin (1977). Maximum likelihood from incomplete data via the EM algorithm. *Journal of the Royal Statistical Society. Series B (Methodological)* 39(1), 1–38.
- Evans, K., T. Love, and S. Thurston (2015). Outlier identification in model-based cluster analysis. *Journal of Classification* 32(1), 63–84.
- Farcomeni, A. and L. Greco (2015). S-estimation of hidden Markov models. *Computational Statistics* 30(1), 57–80.
- Fraley, C. and A. E. Raftery (2002). Model-based clustering, discriminant analysis, and density estimation. *Journal of the American Statistical Association* 97(458), 611–631.
- Fraley, C., A. E. Raftery, L. Scrucca, T. B. Murphy, and M. Fop (2015). *mclust: Normal Mixture Modelling for Model-Based Clustering, Classification, and Density Estimation*. R package version 5.0.1.

- Fritz, H., L. A. García-Escudero, and A. Mayo-Iscar (2012, 5). **tclust**: an R package for a trimming approach to cluster analysis. *Journal of Statistical Software* 47(12), 1–26.
- García-Escudero, L. A., A. Gordaliza, C. Matran, and A. Mayo-Iscar (2015). Avoiding spurious local maximizers in mixture modeling. *Statistics and Computing* 25(3), 619–633.
- Greselin, F., S. Ingrassia, and A. Punzo (2011). Assessing the pattern of covariance matrices via an augmentation multiple testing procedure. *Statistical Methods & Applications* 20(2), 141–170.
- Greselin, F. and A. Punzo (2013). Closed likelihood ratio testing procedures to assess similarity of covariance matrices. *The American Statistician* 67(3), 117–128.
- Hennig, C. (2002). Fixed point clusters for linear regression: computation and comparison. *Journal of Classification* 19(2), 249–276.
- Holzmann, H. and F. Schwaiger (2015). Hidden Markov models with state-dependent mixtures: minimal representation, model testing and applications to clustering. *Statistics and Computing*. DOI: 10.1007/s11222-014-9481-1.
- Humburg, P., D. Bulger, and G. Stone (2008). Parameter estimation for robust HMM analysis of chip-chip data. *BMC Bioinformatics* 9(1), 343.
- Lagona, F., A. Maruotti, and F. Padovano (2015). Multilevel multivariate modelling of legislative count data, with a hidden Markov chain. *Journal of the Royal Statistical Society: Series A (Statistics in Society)* 178, 705–723.
- Leroux, B. G. (1992). Maximum-likelihood estimation for hidden Markov models. *Stochastic processes and their applications* 40(1), 127–143.
- Little, R. J. A. (1988). Robust estimation of the mean and covariance matrix from data with missing values. *Applied Statistics* 37(1), 23–38.
- Maruotti, A. (2011). Mixed hidden Markov models for longitudinal data: an overview. *International Statistical Review* 79(3), 427–454.
- Maruotti, A. (2014). Robust fitting of hidden Markov regression models under a longitudinal setting. *Journal of Statistical Computation and Simulation* 84(8), 1728–1747.

- Maruotti, A. and R. Rocci (2012). A mixed non-homogeneous hidden Markov model for categorical data, with application to alcohol consumption. *Statistics in Medicine* 31(9), 871–886.
- McNicholas, P. D. and T. B. Murphy (2010). Model-based clustering of longitudinal data. *The Canadian Journal of Statistics* 38(1), 153–168.
- Meng, X.-L. and D. B. Rubin (1993). Maximum likelihood estimation via the ECM algorithm: A general framework. *Biometrika* 80(2), 267–278.
- Peel, D. and G. J. McLachlan (2000). Robust mixture modelling using the t distribution. *Statistics and Computing* 10(4), 339–348.
- Punzo, A. and P. D. McNicholas (2014a). Robust clustering in regression analysis via the contaminated Gaussian cluster-weighted model. arXiv.org e-print 1409.6019, available at: <http://arxiv.org/abs/1409.6019>.
- Punzo, A. and P. D. McNicholas (2014b). Robust high-dimensional modeling with the contaminated Gaussian distribution. arXiv.org e-print 1408.2128, available at: <http://arxiv.org/abs/1408.2128>.
- Punzo, A. and P. D. McNicholas (2015). Parsimonious mixtures of contaminated Gaussian distributions with application to allometric studies. arXiv.org e-print 1305.4669, available at: <http://arxiv.org/abs/1305.4669>.
- R Core Team (2013). *R: A Language and Environment for Statistical Computing*. Vienna, Austria: R Foundation for Statistical Computing.
- Rousseeuw, P. J. and A. M. Leroy (2005). *Robust Regression and Outlier Detection*. Wiley.
- Tukey, J. W. (1960). A survey of sampling from contaminated distributions. In I. Olkin (Ed.), *Contributions to Probability and Statistics: Essays in Honor of Harold Hotelling*, Stanford Studies in Mathematics and Statistics, Chapter 39, pp. 448–485. California: Stanford University Press.
- Viroli, C. (2011). Model based clustering for three-way data structures. *Bayesian Analysis* 6(4), 573–602.
- Volant, S., C. Berard, M. L. Martin-Magniette, and S. Robin (2014). Hidden Markov models with mixtures as emission distributions. *Statistics and Computing* 24(4), 493–504.

ACCEPTED MANUSCRIPT

Yao, W. (2012). Model based labeling for mixture models. *Statistics and Computing* 22(2), 337–347.

Zucchini, W. and I. L. MacDonald (2009). *Hidden Markov models for time series: An introduction using R*. Boca Raton, FL: Chapman & Hall.

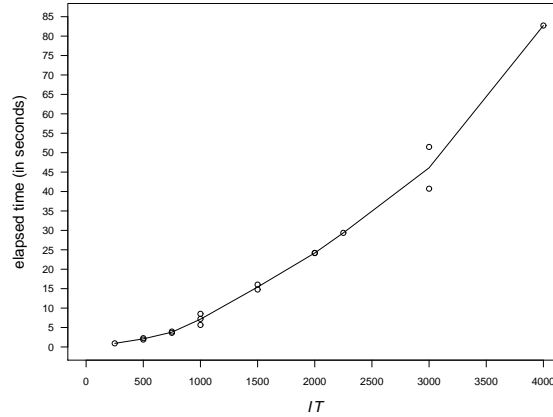


Figure 1: Average elapsed time (in seconds over 100 replications) to fit a CNHMM as a function of IT .

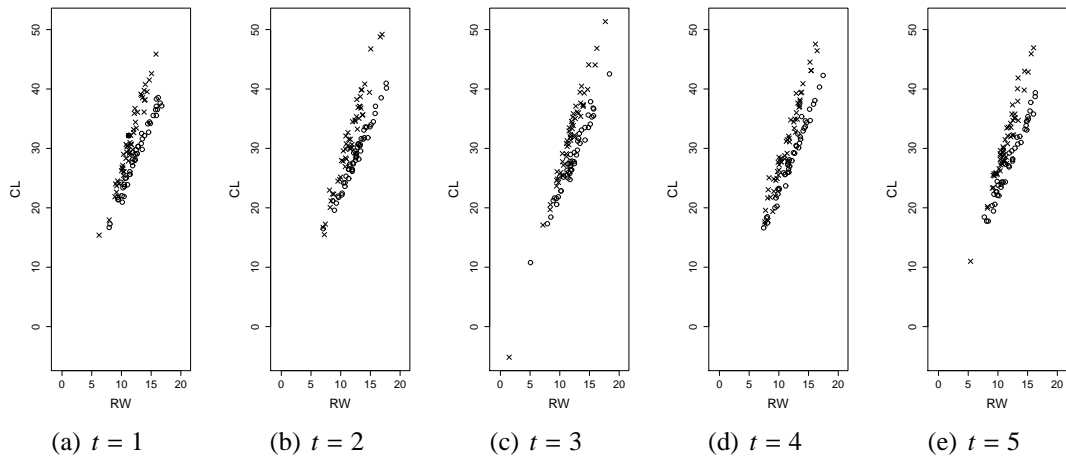


Figure 2: Scatterplots of the artificial data (\times and \circ denote group 1 and group 2, respectively; \bullet denotes the observation perturbed for the analysis of Section 5.1).

	I	T	TPR	FPR		I	T	TPR	FPR
Scenario d)	50	5	1.000	0.003	Scenario e)	50	5	0.860	0.003
		10	1.000	0.002			10	0.841	0.003
		20	1.000	0.000			20	0.844	0.002
	100	5	1.000	0.001		100	5	0.838	0.002
		10	1.000	0.000			10	0.844	0.002
		20	1.000	0.000			20	0.835	0.002
	200	5	1.000	0.000		200	5	0.839	0.002
		10	1.000	0.000			10	0.845	0.002
		20	1.000	0.000			20	0.839	0.002

Table 1: Values of TPRs and FPRs; they refers to rates across 100 replications.

$I \backslash T$	5	10	15	20
50	0.918	1.897	3.955	7.215
100	2.252	8.522	16.039	24.153
150	3.627	14.779	29.357	40.727
200	5.662	24.205	51.474	82.731

Table 2: Average elapsed time (in seconds over 100 replications) to fit a CNHMM as a function of I and T .

Value	NHMM	<i>t</i> HMM	CNHMM	$\hat{\eta}$	\hat{p}
-15	13	5	3	220.966	0.005
-10	12	5	3	174.162	0.006
-5	12	5	3	133.014	0.008
0	11	5	3	97.475	0.010
5	8	5	3	67.539	0.015
10	7	5	3	43.231	0.023
15	7	5	3	24.608	0.041
20	5	5	3	11.729	0.085

Table 3: Number of misallocated artificial blue crabs ($N = 100$ and $T = 5$). The last two columns report the estimated value of η in the group containing the outlier and its posterior probability to be a good point, respectively.

		Number of hidden states (K)									
		1	2	3	4	5	6	7	8	9	10
AIC	CNHMM	-20495.69	-19649.58	-19093.36	-18737.51	-18584.11	-18443.71	-18341.15	-18298.57	-18222.24	-18237.17
	NHMM	-21507.87	-19768.81	-19106.63	-18792.78	-18597.50	-18471.70	-18402.28	-18322.74	-18251.39	-18173.76
	t HMM	-20442.14	-19647.40	-19088.17	-18750.09	-18595.55	-18498.36	-18375.11	-18274.80	-18200.13	-18129.92
BIC	CNHMM	-20537.84	-19741.79	-19240.90	-18945.65	-18858.12	-18788.86	-18762.71	-18801.81	-18812.42	-18919.56
	NHMM	-21544.75	-19850.49	-19238.37	-18979.85	-18845.17	-18785.23	-18786.95	-18783.81	-18794.01	-18803.46
	t HMM	-20481.67	-19734.35	-19227.81	-18947.70	-18856.39	-18827.70	-18778.22	-18756.95	-18766.60	-18785.97
ICL	CNHMM	-20537.84	-19743.67	-19247.93	-18949.27	-18866.21	-18794.69	-18768.71	-18811.05	-18819.57	-18925.40
	NHMM	-21544.75	-19853.99	-19244.35	-18986.99	-18856.12	-18791.33	-18798.01	-18795.69	-18797.59	-18811.22
	t HMM	-20481.67	-19735.65	-19232.32	-18953.05	-18864.92	-18833.48	-18785.28	-18770.45	-18775.57	-18793.86

Table 4: Model selection.

		Hidden states						
		1	2	3	4	5	6	7
μ	HOME	308.965	224.503	152.547	301.355	281.040	146.659	293.847
	TEEN	21.814	15.317	9.545	8.707	22.807	8.672	13.838
	ROB	25.145	27.223	22.736	39.350	94.426	81.336	47.521
	PICK	110.239	160.512	49.022	100.677	651.558	130.300	261.593
α		0.999	0.965	0.999	0.999	0.832	0.775	0.951
η		1.001	5.266	1.014	1.001	2.421	12.370	3.575

Table 5: State-specific parameters.