



# Quantile hidden semi-Markov models for multivariate time series

Luca Merlo<sup>1</sup> · Antonello Maruotti<sup>2,3</sup> · Lea Petrella<sup>4</sup> · Antonio Punzo<sup>5</sup>

Received: 21 September 2021 / Accepted: 24 July 2022 / Published online: 9 August 2022  
© The Author(s), under exclusive licence to Springer Science+Business Media, LLC, part of Springer Nature 2022

## Abstract

This paper develops a quantile hidden semi-Markov regression to jointly estimate multiple quantiles for the analysis of multivariate time series. The approach is based upon the Multivariate Asymmetric Laplace (MAL) distribution, which allows to model the quantiles of all univariate conditional distributions of a multivariate response simultaneously, incorporating the correlation structure among the outcomes. Unobserved serial heterogeneity across observations is modeled by introducing regime-dependent parameters that evolve according to a latent finite-state semi-Markov chain. Exploiting the hierarchical representation of the MAL, inference is carried out using an efficient Expectation-Maximization algorithm based on closed form updates for all model parameters, without parametric assumptions about the states' sojourn distributions. The validity of the proposed methodology is analyzed both by a simulation study and through the empirical analysis of air pollutant concentrations in a small Italian city.

**Keywords** EM algorithm · Latent process · Maximum likelihood · Multivariate asymmetric Laplace distribution · Quantile regression · Sojourn distribution

## 1 Introduction

Since their introduction in the 1960s, Hidden Markov Models (HMMs, see MacDonald and Zucchini 1997; Cappé et al. 2006; Zucchini et al. 2016) have been successfully implemented in a wide range of applications for the analysis of time series data. This class of models is described by an observable stochastic process whose dynamic is governed, completely or partially, by a latent unobservable Markov chain. Owing to their mathematical tractability and the availability of efficient computational procedures, the use of HMMs is well justified when the researcher is interested in inference and/or predictions about the latent process based on the observed

one. For a detailed survey of the literature and fields of application, please see MacDonald and Zucchini 1997; Ephraim and Merhav 2002; Cappé et al. 2006; Maruotti 2011; Bartolucci et al. 2012; Zucchini et al. 2016 and Maruotti and Punzo (2021).

One immediate consequence of the Markov property is that in any HMM, the sojourn time (also defined as state duration or dwell-time), that is, the number of consecutive time points that the Markov chain spends in a given state, is implicitly geometrically distributed (Langrock and Zucchini 2011; Zucchini et al. 2016). Despite the popularity of HMMs, this assumption may not be realistic in many applications which can lead to biased parameter estimates and deteriorate the states classification performance due to a misspecification of the dynamic of the hidden process. Bulla and Bulla (2006) and Maruotti et al. (2019), for instance, show the inability of HMMs to model temporal dependence and reproduce empirical characteristics in real-world data, especially when the probability mass function of sojourn times is far from being geometric.

Motivated by these considerations, Hidden Semi-Markov Models (HSMMs, see Yu 2015) are designed to relax this condition by allowing the Sojourn Distributions (SDs) to be modeled directly by the researcher using more flexible parametric or nonparametric distributions.

✉ Luca Merlo  
luca.merlo@uniroma1.it

<sup>1</sup> Department of Statistical Sciences, Sapienza University of Rome, Piazzale Aldo Moro, 5, 00185 Rome, Italy

<sup>2</sup> Department of Mathematics, University of Bergen, Bergen, Norway

<sup>3</sup> Department of Law, Economics, Political Sciences and Modern Languages, LUMSA University, Rome, Italy

<sup>4</sup> MEMOTEF Department, Sapienza University of Rome, Rome, Italy

<sup>5</sup> Department of Economics and Business, University of Catania, Catania, Italy

Typical choices include families of discrete (semi-)parametric distributions or, alternatively, one can avoid distributional assumptions by estimating the SD probability mass functions based on the observations (Sansom and Thomson 2001; Guédon 2003; Pohle et al. 2022). Combining the increased flexibility to capture a wide range of distributional shapes of the SDs with the well-known advantages of HMMs, HSMMs constitute a versatile framework in several spheres of application (see Guédon 2003; Barbu and Limnios 2009; Bulla et al. 2010; O’Connell and Højsgaard 2011; Yu 2015; Maruotti and Punzo 2021 and the references therein).

In the regression context where covariates are available, both HMMs and HSMMs have also been extended to include a set of predictors by introducing state-dependent regression parameters that evolve over time according to the unobserved process (Hamilton 1989; Yu 2015; Zucchini et al. 2016).

This model specification permits to investigate the dynamics of the hidden state sequence and, at the same time, allows to examine state-specific covariate effects in the observable process, providing a useful modeling framework to capture unobserved time-dependent heterogeneity. Typically, a linear model targeting the conditional mean of the dependent variable given covariates is specified. The assumptions underlying traditional linear regression models, however, are seldom satisfied in real data which often exhibit skewness, heavy tails and outliers. Moreover, the effect of the covariates can differ greatly between different parts of the response distribution. Therefore, when the aim of the research focuses not only at the center of the response distribution but also, and especially in the tails, quantile regression (Koenker and Bassett 1978) represents an interesting alternative to standard mean regression. This method provides a way to model the conditional quantiles of a response variable with respect to the covariates in order to have a more complete picture of the entire conditional distribution compared to ordinary least squares. In the univariate quantile regression framework, both the classical and Bayesian inferential approaches have been proposed in the literature to estimate the model parameters. In the frequentist setting, the inferential approach relies on the minimization of the asymmetric loss function (see Koenker and Bassett 1978) while, in the Bayesian setting and in a likelihood inferential approach, the Asymmetric Laplace (AL) distribution has been introduced as a likelihood inferential tool. The two approaches are well-justified by the relationship between the quantile loss function and the AL density. Indeed, Yu and Moyeed (2001) showed that the minimization of the quantile loss function is equivalent, in terms of parameter estimates, to the maximization of the likelihood associated with the AL density. For a detailed review and list of references, Koenker (2005); Luo et al. (2012); Bernardi et al. (2015) and Koenker et al. (2017) provide an overview of the most used quantile regression techniques in both the classical and Bayesian settings. In addition, quan-

tile regression methods have also been generalized to account for serial heterogeneity. In the analysis of longitudinal data, Farcomeni (2012) and Marino et al. (2018) consider univariate linear quantile models where unobserved sources of time-varying heterogeneity are captured by means of state-dependent coefficients evolving according to a finite-state homogeneous hidden Markov chain. Further, Ye et al. (2016) and Maruotti et al. (2021) propose a (semi-)Markov quantile regression to model the regime-switching effect of the regression coefficients in financial and environmental time series.

When multivariate response variables are concerned, the existing literature on quantile regression is less extensive since there is no “natural” ordering in a  $p$ -dimensional space, for  $p > 1$ . As a consequence, the univariate quantile regression method does not straightforwardly extend to higher dimensions. Nevertheless, in most situations of practical interest, the purpose of the matter being investigated lies in describing the distribution of a multivariate response variable.

For this reason the search for a satisfactory notion of multivariate quantile has led to a flourishing literature on this topic despite its definition is still a debatable issue (see Serfling 2002; Kong and Mizera 2012; Koenker et al. 2017; Stolfi et al. 2018; Chavas 2018; Charlier et al. 2020; Merlo et al. 2021, 2022 and the references therein for relevant studies).

Recently, Petrella and Raponi (2019) generalized the AL distribution inferential approach of the univariate quantile regression to a multivariate framework by using the Multivariate Asymmetric Laplace (MAL) distribution defined in Kotz et al. (2012). Employing the MAL distribution as a likelihood based inferential tool, the authors sidestep the problem of defining the quantiles of a multivariate distribution, and instead implement joint estimation for the univariate quantiles of the conditional distribution of a multivariate response variable given covariates, accounting for possible correlation among the responses.

The purpose of this article is to extend the work of Petrella and Raponi (2019) by introducing a HSMM for the analysis of multivariate time series. More formally, we develop a Quantile Hidden Semi-Markov Model (QHSMM) to jointly estimate the quantiles of the univariate conditional distributions of a multivariate response, accounting for the dependence structure between the outcomes. In particular, to capture the temporal evolution of unobserved heterogeneity, we introduce state-dependent coefficients in the regression model that evolve over time according to a latent semi-Markov process. In order to prevent inconsistent parameter estimates due to misspecification of the SDs, we adopt the nonparametric approach of Guédon (2003) where they are left unspecified and approximated by discrete distributions concentrated on a finite set of time points estimated from the data. Within this scheme, our modeling framework can be thought of as a model-based clustering approach for data

showing time-varying heterogeneity, where the interest lies in the effect of cluster-specific covariates on various quantile levels.

Throughout the paper we propose to estimate the model parameters with a Maximum Likelihood (ML) approach by using the MAL distribution as working likelihood in a regression framework. Specifically, as in Petrella and Raponi (2019) and Merlo et al. (2022), we consider the mixture representation of the MAL distribution which allows us to build an efficient Expectation-Maximization (EM) algorithm with the E- and M-step updates in closed form for all model parameters.

Using simulation experiments, we illustrate the validity of our approach under different data generating processes and evaluate its ability in recovering the true values of the regression coefficients, the true classification and number of latent states.

In the empirical analysis, we apply the proposed methodology to investigate the effect of a collection of atmospheric variables on the daily concentrations of three major pollutants, i.e., particulate matter, ozone and nitrogen dioxide measured in Rieti (Italy) from 2019 to 2021. Our method allows us to: (i) assess how the effects of atmospheric variables can vary across different (more extreme) quantiles of the conditional distribution of air pollutants, accounting for their dependence structure; (ii) summarize the data by means of a reduced number of latent regimes associated with different concentration levels of chemicals.

The paper is organized as follows. In Sect. 2, we introduce the proposed model. Sect. 3 illustrates the EM-based ML approach to estimating the model parameters and the computational details of the algorithm. In Sect. 4 we present the simulation results, while Sect. 5 discusses the empirical application. Finally, Sect. 6 concludes.

## 2 Methodology

Let  $\{S_t\}_{t=1}^T$  denote a finite-state hidden semi-Markov chain defined over a discrete state space  $\mathcal{S} = \{1, \dots, K\}$ . The latent process  $\{S_t\}_{t=1}^T$  is constructed as follows. A homogeneous hidden Markov chain with  $K$  states models the transitions between different states, with initial probabilities,  $\pi_k = \Pr(S_1 = k)$ , and transition probabilities

$$\pi_{jk} = \Pr(S_{t+1} = k \mid S_t = j, S_{t+1} \neq j), \tag{1}$$

with  $\sum_{k=1}^K \pi_{jk} = 1$ ,  $\pi_{jk} \geq 0$ , for every  $k = 1, \dots, K$  and  $\pi_{jj} = 0$ , i.e., the diagonal elements of the transition probability matrix are zeros. More concisely, we collect the initial and transition probabilities in the  $K$ -dimensional vector  $\boldsymbol{\pi}$  and in the  $K \times K$  matrix  $\mathbf{Q}$ , respectively. Because the unobserved process is semi-Markovian, only transitions from one

state to another are governed by the transition probabilities in (1), but the duration of a stay in a state is modeled by a separate SD. Specifically, let us denote by  $d_k(u)$  the SD, i.e., the probability the hidden process  $\{S_t\}_{t=1}^T$  spends  $u$  consecutive time steps in the  $k$ -th state, as follows:

$$d_k(u) = \Pr(S_{t+u} \neq k, S_{t+u-1} = k, \dots, S_{t+1} = k \mid S_t = k, S_{t-1} \neq k), \quad u = 1, \dots, U_k, \tag{2}$$

where  $U_k$  corresponds to the maximum sojourn time of the hidden chain in state  $k$ . Let us also denote  $\mathbf{U} = (U_1, \dots, U_K)$  the  $K$ -dimensional vector collecting all state-specific maximum sojourn times.

HSMMs allow for great flexibility as the SD in (2) is directly specified by the researcher and estimated from the observed data. The SD can be chosen from a large variety of parametric distributions, such as the shifted-Poisson, the shifted-negative binomial distributions or, in the particular case where  $d_k(u)$  is assumed to be geometrically distributed, a HSMM reduces to a HMM with the most likely sojourn time for every state being 1 (Zucchini et al. 2016). Parametric distributions, however, might lack the flexibility to capture key features of empirical SDs in the data, increasing state misclassification rates and inducing substantial bias. Alternatively, semi- and nonparametric data-driven approaches can be adopted (see Sansom and Thomson 2001; Guédon 2003; Langrock and Zucchini 2011; Maruotti et al. 2021) to provide sufficient additional flexibility in comparison to HMMs and accommodate complex distributional shapes.

To build the proposed model, let  $\mathbf{Y}_t = (Y_t^{(1)}, \dots, Y_t^{(p)})'$  be a continuous observable  $p$ -variate response variable and  $\mathbf{X}_t = (1, X_t^{(2)}, \dots, X_t^{(m)})'$  be a  $m$ -dimensional vector of covariates, with the first element being the intercept, at time  $t = 1, \dots, T$ . The process  $\{\mathbf{Y}_t\}_{t=1}^T$  represents the state-dependent process of the HSMM and, conditional on the hidden states, fulfills the independence property:

$$f_{\mathbf{Y}}(\mathbf{y}_t \mid \mathbf{x}_t, \mathbf{y}_1, \dots, \mathbf{y}_{t-1}, S_1 = s_1, \dots, S_t = s_t) = f_{\mathbf{Y}}(\mathbf{y}_t \mid \mathbf{x}_t, S_t = s_t), \tag{3}$$

where  $f_{\mathbf{Y}}(\mathbf{y}_t \mid \mathbf{x}_t, S_t = s_t)$  is the conditional distribution of  $\mathbf{Y}_t$  given the covariates  $\mathbf{X}_t$  and the hidden state occupied at time  $t$ .

As mentioned in Section 1, our objective is to provide joint estimation of the  $p$  quantiles of the univariate conditional distributions of  $\mathbf{Y}_t$ , taking into account time-dependent heterogeneity and potential correlation among the components of  $\mathbf{Y}_t$ . Given  $p$  quantile indexes  $\boldsymbol{\tau} = (\tau_1, \dots, \tau_p)$ , with  $\tau_j \in (0, 1)$ ,  $j = 1, \dots, p$ , the Quantile Hidden Semi-Markov Model (QHSM) is defined as follows:

$$\mathbf{Y}_t = \boldsymbol{\beta}_k(\boldsymbol{\tau})\mathbf{X}_t + \boldsymbol{\epsilon}_{tk}(\boldsymbol{\tau}),$$

$$t = 1, \dots, T \quad \text{and} \quad k = 1, \dots, K, \tag{4}$$

with  $\beta_k(\tau) = (\beta_k^{(1)}(\tau), \dots, \beta_k^{(p)}(\tau))$  being a state-specific  $p \times m$  matrix of unknown regression coefficients that evolves over time according to the hidden process  $S_t$  and takes one of the values in the set  $\{\beta_1(\tau), \dots, \beta_K(\tau)\}$ , and where  $\epsilon_{tk}(\tau)$  denotes a  $p$ -dimensional vector of error terms with univariate component-wise quantiles (at fixed levels  $\tau_1, \dots, \tau_p$ , respectively) equal to zero.

Generalizing the approach of Petrella and Raponi (2019), as conditional distribution of  $\mathbf{Y}_t$  we consider a Multivariate Asymmetric Laplace (MAL) distribution (see Kotz et al. 2012). In detail, based on (4) we assume  $\mathbf{Y}_t | \mathbf{X}_t = \mathbf{x}_t, S_t = k \sim \mathcal{MAL}_p(\mu_{tk}, \mathbf{D}_k \tilde{\xi}, \mathbf{D}_k \Sigma_k \mathbf{D}_k)$  whose probability density function is given by:

$$f_{\mathbf{Y}}(\mathbf{y}_t | \mathbf{x}_t, S_t = k) = \frac{2 \exp \left\{ (\mathbf{y}_t - \mu_{tk})' \mathbf{D}_k^{-1} \Sigma_k^{-1} \tilde{\xi} \right\}}{(2\pi)^{p/2} |\mathbf{D}_k \Sigma_k \mathbf{D}_k|^{1/2}} \left( \frac{\tilde{m}_{tk}}{2 + \tilde{d}_k} \right)^{v/2} K_v \left( \sqrt{(2 + \tilde{d}_k) \tilde{m}_{tk}} \right), \tag{5}$$

where, for each time occasion  $t = 1, \dots, T$  and state  $k = 1, \dots, K$ , the location parameter  $\mu_{tk}$  is defined by the linear model:

$$\mu_{tk} = \mu(S_t = k, \tau) = \beta_k(\tau) \mathbf{X}_t, \tag{6}$$

with  $\mathbf{D}_k \tilde{\xi}$  being the skewness parameter,  $\mathbf{D}_k = \text{diag}[\delta_{k1}, \dots, \delta_{kp}]$ ,  $\delta_{kj} > 0$ , and  $\tilde{\xi} = (\tilde{\xi}_1, \dots, \tilde{\xi}_p)'$  having generic element  $\tilde{\xi}_j = \frac{1-2\tau_j}{\tau_j(1-\tau_j)}$ ,  $j = 1, \dots, p$ .  $\Sigma_k$  is a  $p \times p$  positive definite matrix such that  $\Sigma_k = \Lambda \Psi_k \Lambda$ , with  $\Psi_k$  being a  $p \times p$  state-specific correlation matrix and  $\Lambda = \text{diag}[\sigma_1, \dots, \sigma_p]$ , with  $\sigma_j^2 = \frac{2}{\tau_j(1-\tau_j)}$ ,  $j = 1, \dots, p$ . Moreover,  $\tilde{m}_{tk} = (\mathbf{y}_t - \mu_{tk})' (\mathbf{D}_k \Sigma_k \mathbf{D}_k)^{-1} (\mathbf{y}_t - \mu_{tk})$ ,  $\tilde{d}_k = \mathbf{Q} \Sigma_k^{-1} \tilde{\xi}$ , and  $K_v(\cdot)$  denotes the modified Bessel function of the third kind with index parameter  $v = (2 - p)/2$ .

One of the key benefits of the MAL distribution is that, using (4) and (5), and following Kotz et al. (2012), the  $\mathcal{MAL}_p(\mu, \mathbf{D}\tilde{\xi}, \mathbf{D}\Sigma\mathbf{D})$  can be written as a location-scale mixture with the following representation:

$$\mathbf{Y} = \mu + \mathbf{D}\tilde{\xi}\tilde{C} + \sqrt{\tilde{C}}\mathbf{D}\mathbf{6}^{1/2}\mathbf{Z} \tag{7}$$

where  $\mathbf{Z} \sim \mathcal{N}_p(\mathbf{0}_p, \mathbf{I}_p)$  denotes a  $p$ -variate standard Normal distribution and  $\tilde{C} \sim \text{Exp}(1)$  has a standard exponential distribution, with  $\mathbf{Z}$  being independent of  $\tilde{C}$ . In particular, the constraints imposed on  $\tilde{\xi}$  and  $\Lambda$  guarantee that the  $j$ -th element of  $\mu_{tk}, \mu_{tkj}$ , is the  $\tau_j$ -th conditional quantile function of  $Y_t^{(j)}$  given  $S_t = k$ , for  $k = 1, \dots, K$  and  $j = 1, \dots, p$ , and represent necessary conditions for model identifiability for any fixed quantile level  $\tau_1, \dots, \tau_p$ , as stated in the next proposition.

**Proposition 1** Let  $\mathbf{Y}_t | \mathbf{X}_t = \mathbf{x}_t, S_t = k \sim \mathcal{MAL}_p(\mu_{tk}, \mathbf{D}_k \tilde{\xi}, \mathbf{D}_k \Sigma_k \mathbf{D}_k)$ , for  $k = 1, \dots, K$ , where  $K$  is a positive integer,  $\tilde{\xi} = (\tilde{\xi}_1, \dots, \tilde{\xi}_p)'$  with  $\tilde{\xi}_j = \frac{1-2\tau_j}{\tau_j(1-\tau_j)}$  being known for any fixed value of  $\tau_j$ . Furthermore,  $\Sigma_k = \Lambda \Psi_k \Lambda$  is a  $p \times p$  positive definite matrix with  $\Psi_k$  being an unknown  $p \times p$  correlation matrix and  $\Lambda = \text{diag}[\sigma_1, \dots, \sigma_p]$ , with fixed element  $\sigma_j^2 = \frac{2}{\tau_j(1-\tau_j)}$ ,  $j = 1, \dots, p$ . Then, the model in (2)-(6) is identified.

**Proof** See Proof of Proposition 1 in Appendix. □

In comparison with other methods in the literature, the proposed modeling framework includes the homogeneous

joint quantile regression approach of Petrella and Raponi (2019) when  $K = 1$  and it reduces to the univariate hidden semi-Markov-switching quantile regression of Maruotti et al. (2021) when  $p = 1$ . Naturally, when a geometric SD is assumed for all latent states, we call our methodology the Quantile Hidden Markov Model (QHMM).

### 3 Maximum likelihood estimation and inference

In this section we introduce a ML approach to making inference on model parameters. As is usually done in the literature in the presence of latent variables, we propose a suitable likelihood-based EM algorithm (Dempster et al. 1977). To hedge against possibly biased inference from incorrect parametric assumptions on the SDs, we estimate the sojourn probabilities nonparametrically following Guédon (2003). In addition, we show that both the E- and M-step updates of the algorithm can be obtained in closed form by exploiting the hierarchical representation of the MAL distribution in (7) under the constraints on  $\mathbf{Q}$  and  $\Lambda$ , hence reducing the computational burden compared to direct maximization of the likelihood. We illustrate the EM algorithm to fit the more general QHSMM but it can also be employed for the QHMM by assuming a geometric SD. To ease the notation, unless specified otherwise, we omit the quantile levels vector  $\tau$ , yet all model parameters are allowed to depend on it. All the proofs are collected in the Appendix.

Let us denote by  $\Phi_\tau = (\beta_1, \dots, \beta_K, \mathbf{D}_1, \dots, \mathbf{D}_K, \Psi_1, \dots, \Psi_K, \boldsymbol{\pi}, \mathbf{Q}, d_1(u), \dots, d_K(U_K))$  the set of model parameters. For any fixed  $\tau$ , number of hidden states  $K$  and maximum sojourn times  $\mathbf{U}$ , we use the MAL representation



in (7) to express the complete-data likelihood as follows:

$$L_c(\Phi_\tau) = \pi_{s_1^*} d_{s_1^*}(u_1) \left\{ \prod_{r=2}^{R-1} \pi_{s_r^* | s_{r-1}^*} d_{s_r^*}(u_r) \right\} \pi_{s_R^* | s_{R-1}^*} D_{s_R^*}(u_R) \prod_{t=1}^T f_Y(\mathbf{y}_t | \mathbf{x}_t, s_t, \tilde{c}_t, \tau) f_{\tilde{C}}(\tilde{c}_t), \tag{8}$$

where  $\tilde{C}$  is a latent variable that follows an exponential distribution with parameter 1,  $s_r^*$  is the  $r$ -th visited state,  $u_r$  is the time spent in that state (i.e., the duration of the  $r$ -th visit) and  $R - 1$  is the number of state changes up to time  $T$ . Following Guédon (2003), the survivor function  $D_k(u)$  for the sojourn time in state  $k$  is defined as:

$$D_k(u) = \sum_{v \geq u} d_k(v). \tag{9}$$

The survivor function sums up the individual probability masses of all possible sojourns of length  $v \geq u$  and it has

$$\begin{aligned} \mathcal{O}(\Phi_\tau) = & \sum_{k=1}^K \gamma_{1k} \log \pi_k + \sum_{t=1}^T \sum_{j=1}^K \sum_{k \neq j}^K v_{tjk} \log \pi_{jk} + \sum_{k=1}^K \sum_{u=1}^{U_k} \eta_k(u) \log d_k(u) - \frac{1}{2} T \sum_{k=1}^K \log |\mathbf{D}_k \Sigma_k \mathbf{D}_k| \\ & + \sum_{t=1}^T \sum_{k=1}^K \gamma_{tk} \tilde{z}_{tk} (\mathbf{Y}_t - \boldsymbol{\mu}_{tk})' \mathbf{D}_k^{-1} \Sigma_k^{-1} \tilde{\boldsymbol{\xi}} - \frac{1}{2} \sum_{t=1}^T \sum_{k=1}^K \gamma_{tk} \tilde{z}_{tk} (\mathbf{Y}_t - \boldsymbol{\mu}_{tk})' (\mathbf{D}_k \Sigma_k \mathbf{D}_k)^{-1} (\mathbf{Y}_t - \boldsymbol{\mu}_{tk}) - \frac{1}{2} \sum_{t=1}^T \sum_{k=1}^K \gamma_{tk} \tilde{c}_{tk} \tilde{\boldsymbol{\xi}}' \Sigma_k^{-1} \tilde{\boldsymbol{\xi}}, \end{aligned} \tag{13}$$

several advantages. Firstly, we do not have to assume that the process is leaving a state immediately after the upper endpoint  $T$ . Secondly, it provides a more accurate prediction of the last state visited, which is important when the data analysis wishes to estimate the most recently visited state, and improves parameter estimation (O’Connell and Højsgaard 2011).

### 3.1 The EM algorithm

The EM algorithm alternates between performing an expectation (E) step, which defines the expectation of the complete log-likelihood function evaluated using the current estimates of the parameters, and a maximization (M) step, which computes parameter estimates by maximizing the expected complete log-likelihood obtained in the E-step. The expected complete log-likelihood function and the optimal parameter updates are given in the following propositions.

Given the representation in (8), for the implementation of the algorithm we introduce the following quantities. We define the probability of being in state  $k$  at time  $t$  given the observed sequence as:

$$\gamma_{tk} = \Pr(S_t = k | \mathbf{Y}_1, \dots, \mathbf{Y}_T). \tag{10}$$

The probability the process left state  $j$  at time  $t - 1$  and entered state  $k$  at time  $t$  given the observed sequence is:

$$v_{tjk} = \Pr(S_{t-1} = j, S_t = k | \mathbf{Y}_1, \dots, \mathbf{Y}_T). \tag{11}$$

Finally, let us denote by  $\eta_k(u)$  the expected number of times the process spends  $u$  consecutive time steps in state  $k$  as:

$$\begin{aligned} \eta_k(u) = & \Pr(S_{1+u} \neq k, S_{1+u-v} = k, \\ & v = 1, \dots, u | \mathbf{Y}_1, \dots, \mathbf{Y}_T) \\ & + \sum_{t=2}^T \Pr(S_{t+u} \neq k, S_{t+u-v} = k, v = 1, \dots, u, \\ & S_{t-1} \neq k | \mathbf{Y}_1, \dots, \mathbf{Y}_T). \end{aligned} \tag{12}$$

Then, the expected log-likelihood for the complete data is presented in the following proposition.

**Proposition 2** For any fixed  $\tau = (\tau_1, \dots, \tau_p)$ , number of hidden states  $K$  and maximum sojourn times  $\mathbf{U} = (U_1, \dots, U_K)$ , the expected complete log-likelihood function (up to additive constants) is:

where

$$\begin{aligned} \tilde{c}_{tk} = & \left( \frac{\tilde{m}_{tk}}{2 + \tilde{d}_k} \right)^{\frac{1}{2}} \frac{K_{v+1} \left( \sqrt{(2 + \tilde{d}_k) \tilde{m}_{tk}} \right)}{K_v \left( \sqrt{(2 + \tilde{d}_k) \tilde{m}_{tk}} \right)}, \\ \tilde{z}_{tk} = & \left( \frac{2 + \tilde{d}_k}{\tilde{m}_{tk}} \right)^{\frac{1}{2}} \frac{K_{v+1} \left( \sqrt{(2 + \tilde{d}_k) \tilde{m}_{tk}} \right)}{K_v \left( \sqrt{(2 + \tilde{d}_k) \tilde{m}_{tk}} \right)} - \frac{2v}{\tilde{m}_{tk}}, \end{aligned} \tag{14}$$

with

$$\tilde{m}_{tk} = (\mathbf{y}_t - \boldsymbol{\mu}_{tk})' (\mathbf{D}_k \Sigma_k \mathbf{D}_k)^{-1} (\mathbf{y}_t - \boldsymbol{\mu}_{tk}), \quad \tilde{d}_k = \tilde{\boldsymbol{\xi}}' \Sigma_k^{-1} \tilde{\boldsymbol{\xi}}. \tag{15}$$

Therefore, the EM algorithm can be implemented as follows:

*E-step:* At the generic  $r$ -th iteration of the algorithm, let  $\hat{\Phi}_\tau^{(r-1)}$  denote the current parameter estimates. Then, conditionally on the observed data and  $\hat{\Phi}_\tau^{(r-1)}$ , the quantities,  $\gamma_{tk}$  in (10) and  $v_{tjk}$  in (11), can be calculated via a dynamic programming method known as the forward-backward algorithm (see, e.g., Levinson et al. 1983), while  $\eta_k(u)$  in (12)

can be computed using the efficient adaptation of the forward-backward algorithm provided by Guédon (2003). Similarly, the conditional expectations  $\tilde{c}_{tk}$  and  $\tilde{z}_{tk}$  in (14) are considered; see the Appendix. We denote such quantities as  $\hat{\gamma}_{tk}^{(r)}$ ,  $\hat{v}_{ijk}^{(r)}$ ,  $\hat{\eta}_k^{(r)}(u)$ ,  $\hat{c}_{tk}^{(r)}$  and  $\hat{z}_{tk}^{(r)}$ .

**M-step:** Substitute them in (13) to maximize  $\mathcal{O}(\Phi_\tau | \hat{\Phi}_\tau^{(r-1)})$  with respect to  $\Phi_\tau$ , and obtain the updated parameter estimates. Because the expected complete-data log-likelihood in (13) decomposes into orthogonal subproblems, the maximization with respect to the regression coefficients, the parameters of the MAL distribution and the hidden process, can be performed separately. Thus, the initial probabilities  $\pi_j$  and transition probabilities  $\pi_{jk}$  are estimated by:

$$\hat{\pi}_j^{(r)} = \hat{\gamma}_{1j}^{(r)} \quad \text{and} \quad \hat{\pi}_{jk}^{(r)} = \frac{\sum_{t=1}^T \hat{v}_{tjk}^{(r)}}{\sum_{t=1}^T \sum_{j \neq k}^K \hat{v}_{tjk}^{(r)}}. \tag{16}$$

To update the state-specific SD, we follow the nonparametric approach of Guédon (2003). In particular, we set the latent state duration densities to be discrete nonparametric distributions with arbitrary point mass assigned to the feasible duration values, that is, the SD is estimated as follows:

$$\hat{d}_k^{(r)}(u) = \frac{\hat{\eta}_k^{(r)}(u)}{\sum_{v=1}^{U_k} \hat{\eta}_k^{(r)}(v)}. \tag{17}$$

Finally, the M-step updates of the parameters in the regression equation  $\{\beta_j, \mathbf{D}_j, \Sigma_j\}_{j=1}^K$ , are given in the following proposition.

**Proposition 3** *At the generic  $r$ -th iteration, the values of  $\beta_j, \mathbf{D}_j$  and  $\Sigma_j$  maximizing (13) are:*

$$\hat{\beta}_j^{(r)} = \left( \sum_{t=1}^T \hat{\gamma}_{tj}^{(r)} \hat{z}_{tj}^{(r)} \mathbf{X}_t \mathbf{X}_t' \right)^{-1} \left( \sum_{t=1}^T \hat{\gamma}_{tj}^{(r)} \hat{z}_{tj}^{(r)} \mathbf{X}_t \mathbf{Y}_t' - \sum_{t=1}^T \hat{\gamma}_{tj}^{(r)} \mathbf{X}_{it} \tilde{\xi}' \hat{\mathbf{D}}_j^{(r-1)} \right). \tag{18}$$

$$\hat{\Sigma}_j^{(r)} = \frac{1}{T} \sum_{t=1}^T \hat{\gamma}_{tj}^{(r)} \hat{z}_{tj}^{(r)} \hat{\mathbf{D}}_j^{-1(r-1)} (\mathbf{Y}_t - \hat{\mu}_{tj}^{(r)}) (\mathbf{Y}_t - \hat{\mu}_{tj}^{(r)})' \hat{\mathbf{D}}_j^{-1(r-1)} + \frac{1}{T} \sum_{t=1}^T \hat{\gamma}_{tj}^{(r)} \hat{c}_{tj}^{(r)} \tilde{\xi} \tilde{\xi}' - \frac{2}{T} \hat{\mathbf{D}}_j^{-1(r-1)} \sum_{t=1}^T \hat{\gamma}_{tj}^{(r)} (\mathbf{Y}_t - \hat{\mu}_{tj}^{(r)}) \tilde{\xi}', \tag{19}$$

where  $\hat{\mu}_{tj}^{(r)} = \hat{\beta}_j^{(r)} \mathbf{X}_t$ .

For the  $j$ -th state, the elements  $\delta_{jk}, k = 1, \dots, p$ , of the diagonal scale matrix  $\mathbf{D}_j$  are estimated by:

$$\hat{\delta}_{jk}^{(r)} = \frac{1}{T} \sum_{t=1}^T \hat{\gamma}_{tj}^{(r)} \rho_\tau(Y_t^{(k)} - \hat{\mu}_{tjk}^{(r)}), \tag{20}$$

where  $\rho_\tau(\cdot)$  is the quantile check function of Koenker and Bassett (1978):

$$\rho_\tau(u) = u(\tau - \mathbf{1}(u < 0)), \tag{21}$$

with  $\mathbf{1}(\cdot)$  being the indicator function and  $\hat{\mu}_{tjk}^{(r)}$  being the  $k$ -th element of the vector  $\hat{\mu}_{tj}^{(r)}$ .

The E- and M-steps are alternated until convergence, that is when the observed likelihood between two consecutive iterations is smaller than a predetermined threshold. In this paper, we set this threshold criterion equal to  $10^{-5}$ .

Following Maruotti et al. (2021), for fixed  $\tau, K$  and  $\mathbf{U}$ , we initialize the EM algorithm by providing the initial states partition,  $\{S_t^{(0)}\}_{t=1}^T$ , according to a Multinomial distribution with probabilities  $1/K$ . From the generated partition, the off-diagonal elements of  $\hat{\mathbf{Q}}^{(0)}$  are computed as proportions of transition. We obtain  $\hat{\beta}_k^{(0)}$  and  $\hat{\mathbf{D}}_k^{(0)}$  by fitting univariate quantile regressions on observations within state  $k$ , while  $\hat{\Psi}_k^{(0)}$  is set equal to the empirical correlation computed on observations in the  $k$ -th state. The initial SDs are estimated from  $\{S_t^{(0)}\}_{t=1}^T$  assuming a geometric distribution as in HMMs. To avoid convergence to local maxima and better explore the parameter space, we fit the proposed QHSMM using a multiple random starts strategy with different starting partitions and retain the solution corresponding to the maximum likelihood value.

Once we computed the ML estimates of the model parameters  $\hat{\Phi}_\tau$ , we calculate standard errors using a parametric bootstrap approach (Visser et al. 2000). That is, we refitted the model to  $H$  bootstrap samples and approximate the standard error of each model parameter with its corresponding standard deviation computed on bootstrap samples. Hence, standard error estimates for  $\hat{\Phi}_\tau$  are given by the diagonal elements of:

$$\widehat{\text{Cov}}(\hat{\Phi}_\tau) = \sqrt{\frac{1}{H-1} \sum_{h=1}^H (\hat{\Phi}_\tau^{(h)} - \bar{\Phi}_\tau) (\hat{\Phi}_\tau^{(h)} - \bar{\Phi}_\tau)'}, \tag{22}$$

where  $\hat{\Phi}_\tau^{(h)}$  is the set of parameter estimates for the  $h$ -th bootstrap sample and  $\bar{\Phi}_\tau$  denote the sample mean of all  $\hat{\Phi}_\tau^{(h)}, h = 1, \dots, H$ .

### 3.2 Model selection

In the EM algorithm discussed above, the number of hidden states  $K$  and maximum length of state durations  $\mathbf{U} = (U_1, \dots, U_K)$  are unknown. From an applied perspective, choosing an adequate number of states is a crucial aspect of the data analysis which shall take into account the data structure and research question at hand. In particular,  $K$  is typically selected using penalized-likelihood criteria or cross-validation methods, which can become demanding to fit computationally (Pohle et al. 2017). In HSMs, not only the number of hidden states shall be selected, but also the maximum length of state durations. In practice, this amounts to fixing  $U_k = U, k = 1, \dots, K$ , with  $U$  being large enough to capture the main support of the SD in each state (see Maruotti et al. 2021). The major disadvantages of this approach are the large number of parameters to be estimated and the fact that different states may require substantially different maximum sojourn times. For these reasons, herein we simultaneously select the optimal values of  $K$  and vector  $\mathbf{U}$  using penalized likelihood criteria, such as AIC (Akaike 1998), BIC (Schwarz 1978) and ICL (Biernacki et al. 2000) which penalizes the BIC for the estimated mean entropy and it is given by:

$$ICL_{(K,U)} = BIC_{(K,U)} - 2 \sum_{t=1}^T \sum_{k=1}^K \hat{\gamma}_{tk} \log \hat{\gamma}_{tk}, \tag{23}$$

where  $BIC_{(K,U)}$  in (23) is defined as  $BIC_{(K,U)} = -2\ell(\hat{\Phi}_\tau) + \log(T)v_f$ , with  $\ell(\hat{\Phi}_\tau) = \log \sum_{s_1, \dots, s_T} \sum_u L_c(\hat{\Phi}_\tau)$  being the observed data log-likelihood in correspondence of the ML estimate of  $\Phi_\tau$ ,  $T$  corresponds to the number of observations and  $v_f$  denotes the number of free model parameters in  $\Phi_\tau$ . Computing  $\ell(\hat{\Phi}_\tau)$  may prove to be difficult to evaluate directly because it involves the sum on every possible state sequence of length  $T$ ,  $\sum_{s_1, \dots, s_T}$ , and the sum on every supplementary duration from time  $T + 1$  spent in the state occupied at time  $T$ ,  $\sum_u$ . Therefore, to compute  $\ell(\hat{\Phi}_\tau)$  we use the variables in (10)-(12) required for the EM algorithm (please see Guédon 2003).

All criteria involve penalization terms depending on the number of parameters  $v_f$ , which is given by the sum of:

- the number of regression parameters in  $\{\beta_1, \dots, \beta_K\}$ :  $p \times m \times K$ ,
- the number of scale parameters in  $\mathbf{D}_k$ :  $p \times K$ ,
- the number of correlation parameters in  $\Psi_k$ :  $p \times K(K - 1)/2$ ,
- the number of independent transition probabilities in  $\mathbf{Q}$ :  $K \times (K - 2)$ ,
- the unconstrained sojourn distribution probabilities:  $\sum_{k=1}^K (U_k - 1)$ .

To select the order of the hidden process, we first define a sequence of values of  $K$  and construct a  $K$ -dimensional grid of maximum sojourn distributions,  $\mathcal{U} \subset \mathbb{R}_{\geq 0}^K$ , and then fit the model using the EM algorithm described above for fixed  $\tau$ ,  $K$  and a vector  $\mathbf{U}$  in  $\mathcal{U}$ . Because a full search over  $\mathcal{U}$  might be computationally infeasible, we employ the greedy search algorithm considered in Langrock et al. (2015) and Adam et al. (2019) and select the best combination of  $(K, \mathbf{U})$  corresponding to the lowest value of the penalized likelihood criteria.

### 4 Simulation study

We conduct a simulation study to evaluate the finite sample properties of the proposed QHSM. This simulation exercise addresses the following issues: (i) study the performance of the model under different distributional choices for the error term and SDs, when either a linear or nonlinear quantile regression function of  $\mathbf{Y}$  given  $\mathbf{X}$  is considered; (ii) assess the classification performance of the proposed model; (iii) evaluate the performance of penalized likelihood criteria in selecting the optimal number of hidden states  $K$  and maximum sojourn times  $\mathbf{U}$ . Additional simulation studies are illustrated in the Supplementary Materials.

We consider  $T = 1000$ , a continuous response variable of dimension  $p = 2$  and one explanatory variable  $X_t \sim \mathcal{N}(0, 1)$ . The observations are generated from a two state HSM, i.e.,  $K = 2$ , using the following data generating process:

$$\mathbf{Y}_t = \beta_k \mathbf{X}_t + \epsilon_{tk}, \tag{24}$$

where  $\mathbf{X}_t = (1, X_t)'$  and the true values of the state-dependent parameters,  $\beta_1$  and  $\beta_2$ , are given by:

$$\beta_1 = \begin{pmatrix} 4 & 2 \\ -3 & -1 \end{pmatrix} \quad \text{and} \quad \beta_2 = \begin{pmatrix} 5 & -2 \\ -4 & 1 \end{pmatrix}. \tag{25}$$

We consider the following two distributions for the error terms  $\epsilon_{tk}$  in (24):

- ( $\mathcal{N}$ ) :  $\epsilon_{tk}$  are generated from a bivariate Normal random variable with zero mean vector and variance-covariance matrix equal to  $\tilde{\Omega}_k$ , for  $k = 1, 2$ ;
- ( $\mathcal{T}$ ) :  $\epsilon_{tk}$  are generated from a bivariate Student t distribution with 5 degrees of freedom, zero mean and scale matrix equal to  $\tilde{\Omega}_k$ , for  $k = 1, 2$ .

The state-specific covariance matrices  $\tilde{\Omega}_k, k = 1, 2$ , are set equal to low ( $\tilde{\Omega}_1 = \begin{pmatrix} 1 & 0.3 \\ 0.3 & 1 \end{pmatrix}$ ) and high ( $\tilde{\Omega}_2 = \begin{pmatrix} 1 & 0.8 \\ 0.8 & 1 \end{pmatrix}$ ) correlation between the responses.

**Table 1** ARB (in percentage) and RMSE (in brackets) for state-parameter estimates of  $\beta_1$  and  $\beta_2$  with normal errors for the QHSMM

| True Coef.          | $\tau$         |                |                |                |                |
|---------------------|----------------|----------------|----------------|----------------|----------------|
|                     | (0.10, 0.10)   | (0.25, 0.25)   | (0.50, 0.50)   | (0.75, 0.75)   | (0.90, 0.90)   |
| <i>Panel A: SPO</i> |                |                |                |                |                |
| -4                  | 2.584 (0.081)  | 0.653 (0.062)  | 0.095 (0.058)  | -0.452 (0.066) | -2.227 (0.090) |
| -3                  | 3.576 (0.051)  | 0.447 (0.039)  | -0.053 (0.035) | -0.549 (0.038) | -3.771 (0.053) |
| -2                  | 0.243 (0.081)  | 0.149 (0.068)  | 0.150 (0.065)  | 0.040 (0.070)  | 0.214 (0.086)  |
| -1                  | 0.055 (0.064)  | 0.146 (0.052)  | 0.201 (0.050)  | 0.161 (0.051)  | 0.293 (0.065)  |
| 1                   | -0.588 (0.085) | -0.344 (0.068) | -0.271 (0.067) | -0.074 (0.071) | -0.344 (0.088) |
| 2                   | 0.047 (0.063)  | -0.085 (0.051) | -0.152 (0.050) | -0.205 (0.051) | -0.265 (0.063) |
| 4                   | -2.658 (0.052) | -0.342 (0.038) | 0.023 (0.035)  | 0.374 (0.039)  | 2.674 (0.054)  |
| 5                   | -1.994 (0.081) | -0.487 (0.061) | -0.077 (0.058) | 0.415 (0.062)  | 1.797 (0.085)  |
| <i>Panel B: SNB</i> |                |                |                |                |                |
| -4                  | 2.851 (0.092)  | 0.738 (0.072)  | 0.100 (0.067)  | -0.579 (0.077) | -2.414 (0.110) |
| -3                  | 3.390 (0.050)  | 0.424 (0.036)  | -0.021 (0.032) | -0.493 (0.035) | -3.689 (0.050) |
| -2                  | 0.298 (0.097)  | 0.457 (0.079)  | 0.044 (0.073)  | -0.179 (0.080) | -0.061 (0.098) |
| -1                  | -0.335 (0.061) | 0.111 (0.048)  | 0.087 (0.047)  | -0.190 (0.050) | -0.062 (0.061) |
| 1                   | -0.457 (0.098) | -0.867 (0.078) | 0.007 (0.071)  | 0.352 (0.078)  | 0.257 (0.096)  |
| 2                   | 0.255 (0.061)  | -0.043 (0.050) | -0.131 (0.047) | -0.010 (0.049) | 0.073 (0.063)  |
| 4                   | -2.496 (0.050) | -0.274 (0.036) | 0.047 (0.033)  | 0.323 (0.036)  | 2.545 (0.052)  |
| 5                   | -2.252 (0.095) | -0.578 (0.074) | -0.076 (0.070) | 0.488 (0.074)  | 1.881 (0.103)  |
| <i>Panel C: GEO</i> |                |                |                |                |                |
| -4                  | 2.771 (0.076)  | 0.774 (0.060)  | 0.027 (0.056)  | -0.589 (0.065) | -2.306 (0.090) |
| -3                  | 3.174 (0.062)  | 0.235 (0.044)  | -0.037 (0.040) | -0.556 (0.044) | -3.769 (0.064) |
| -2                  | 0.030 (0.083)  | -0.033 (0.068) | -0.082 (0.063) | -0.068 (0.068) | 0.090 (0.082)  |
| -1                  | -1.221 (0.069) | -0.284 (0.053) | -0.134 (0.049) | -0.055 (0.053) | 0.189 (0.067)  |
| 1                   | -0.292 (0.082) | -0.231 (0.067) | -0.154 (0.063) | -0.189 (0.069) | -0.249 (0.081) |
| 2                   | 0.841 (0.070)  | 0.264 (0.055)  | -0.019 (0.051) | -0.015 (0.052) | -0.005 (0.066) |
| 4                   | -2.420 (0.059) | -0.219 (0.044) | 0.038 (0.039)  | 0.300 (0.044)  | 2.417 (0.061)  |
| 5                   | -2.137 (0.076) | -0.619 (0.060) | -0.041 (0.058) | 0.474 (0.062)  | 1.779 (0.087)  |

Similarly to Maruotti et al. (2021), for each scenario we further consider three SDs:

(SPO) : a shifted-Poisson, i.e.:

$$d_k(u) = \exp(-\lambda_k) \frac{\lambda_k^{u-1}}{(u-1)!}, \quad u = 1, 2, \dots, \quad (26)$$

with  $\lambda_1 = 10$  and  $\lambda_2 = 5$ ;

(SNB) : a shifted-negative binomial, i.e.:

$$d_k(u) = \frac{\Gamma(u + \lambda_k - 1)}{(u-1)! \Gamma(\lambda_k)} p_k^{\lambda_k} (1 - p_k)^{u-1}, \quad u = 1, 2, \dots, \quad (27)$$

with  $\lambda_1 = 8, \lambda_2 = 4, p_1 = 0.5$  and  $p_2 = 0.6$ ;

(GEO) : a geometric sojourn, i.e.:

$$d_k(u) = p_k(1 - p_k)^{u-1}, \quad u = 1, 2, \dots, \quad (28)$$

with  $p_1 = 0.2$  and  $p_2 = 0.3$ .

We fit the proposed QHSMM for five quantile levels, i.e.,  $\tau = (0.10, 0.10), \tau = (0.25, 0.25), \tau = (0.50, 0.50), \tau = (0.75, 0.75)$  and  $\tau = (0.90, 0.90)$ . For each model, we carry out  $B = 1000$  Monte Carlo replications and report the following indicators. The Average Relative Bias (ARB), expressed as a percentage:

$$ARB(\hat{\theta}_\tau) = \frac{1}{B} \sum_{b=1}^B \frac{(\hat{\theta}_\tau^{(b)} - \theta_\tau)}{\theta_\tau} \times 100, \quad (29)$$

where  $\hat{\theta}_\tau^{(b)}$  is the estimated parameter at level  $\tau$  for the  $b$ -th replication and  $\theta_\tau$  is the corresponding “true” value. Secondly, the Root Mean Square Error (RMSE) of model



**Table 2** ARB (in percentage) and RMSE (in brackets) for state-parameter estimates of  $\beta_1$  and  $\beta_2$  with Student t errors for the QHSMM

| True Coef.          | $\tau$         |                |                |                |                |
|---------------------|----------------|----------------|----------------|----------------|----------------|
|                     | (0.10, 0.10)   | (0.25, 0.25)   | (0.50, 0.50)   | (0.75, 0.75)   | (0.90, 0.90)   |
| <i>Panel A: SPO</i> |                |                |                |                |                |
| -4                  | 4.150 (0.132)  | 0.393 (0.068)  | 0.082 (0.059)  | -0.112 (0.069) | -3.457 (0.140) |
| -3                  | 3.213 (0.080)  | -0.420 (0.044) | -0.099 (0.037) | 0.200 (0.044)  | -4.030 (0.081) |
| -2                  | 0.178 (0.119)  | -0.098 (0.077) | -0.071 (0.069) | 0.098 (0.078)  | 0.323 (0.124)  |
| -1                  | -0.470 (0.088) | 0.065 (0.058)  | -0.164 (0.050) | -0.221 (0.057) | -0.082 (0.089) |
| 1                   | -0.610 (0.119) | -0.135 (0.077) | -0.273 (0.070) | -0.469 (0.079) | -0.740 (0.126) |
| 2                   | 0.613 (0.086)  | 0.258 (0.056)  | 0.257 (0.048)  | 0.250 (0.056)  | 0.231 (0.087)  |
| 4                   | -2.558 (0.075) | 0.294 (0.042)  | 0.083 (0.036)  | -0.190 (0.042) | 2.500 (0.075)  |
| 5                   | -2.986 (0.132) | -0.236 (0.071) | -0.075 (0.060) | 0.115 (0.069)  | 2.574 (0.137)  |
| <i>Panel B: SNB</i> |                |                |                |                |                |
| -4                  | 5.410 (0.158)  | 0.579 (0.085)  | -0.054 (0.072) | -0.577 (0.084) | -5.037 (0.199) |
| -3                  | 2.516 (0.080)  | -0.423 (0.042) | 0.075 (0.034)  | 0.357 (0.041)  | -3.445 (0.082) |
| -2                  | 0.906 (0.135)  | 0.177 (0.086)  | -0.101 (0.076) | -0.038 (0.087) | -0.171 (0.143) |
| -1                  | -0.945 (0.084) | -0.119 (0.056) | -0.168 (0.050) | -0.059 (0.055) | 0.448 (0.085)  |
| 1                   | -1.643 (0.137) | -0.275 (0.086) | 0.206 (0.075)  | 0.359 (0.085)  | 0.744 (0.141)  |
| 2                   | 0.840 (0.085)  | 0.066 (0.056)  | -0.127 (0.048) | -0.022 (0.054) | -0.047 (0.084) |
| 4                   | -2.102 (0.075) | 0.291 (0.042)  | -0.037 (0.034) | -0.395 (0.040) | 1.812 (0.076)  |
| 5                   | -3.836 (0.171) | -0.391 (0.087) | 0.015 (0.072)  | 0.444 (0.083)  | 3.787 (0.186)  |
| <i>Panel C: GEO</i> |                |                |                |                |                |
| -4                  | 4.668 (0.129)  | 0.717 (0.071)  | -0.012 (0.059) | -0.421 (0.071) | -3.923 (0.167) |
| -3                  | 2.407 (0.102)  | -0.537 (0.052) | 0.015 (0.043)  | 0.308 (0.050)  | -3.759 (0.105) |
| -2                  | -0.068 (0.115) | -0.039 (0.075) | -0.130 (0.065) | -0.173 (0.074) | 0.375 (0.118)  |
| -1                  | -1.090 (0.092) | -0.022 (0.060) | -0.091 (0.052) | 0.119 (0.059)  | 0.323 (0.098)  |
| 1                   | 0.275 (0.118)  | 0.076 (0.076)  | 0.127 (0.064)  | 0.545 (0.075)  | -0.229 (0.118) |
| 2                   | 1.057 (0.093)  | 0.229 (0.062)  | -0.094 (0.053) | -0.068 (0.062) | 0.034 (0.098)  |
| 4                   | -2.102 (0.092) | 0.429 (0.050)  | 0.045 (0.042)  | -0.368 (0.052) | 2.027 (0.098)  |
| 5                   | -3.366 (0.135) | -0.429 (0.072) | -0.016 (0.061) | 0.287 (0.073)  | 2.963 (0.158)  |

parameters averaged across the  $B$  simulations:

$$RMSE(\hat{\theta}_\tau) = \sqrt{\frac{1}{B} \sum_{b=1}^B (\hat{\theta}_\tau^{(b)} - \theta_\tau)^2}. \tag{30}$$

To assess the first and second queries of this simulation exercise, Tables 1 and 2 report the ARB and RMSE for the state-specific coefficients  $\beta_1$  and  $\beta_2$ . As can be noted, the proposed model under the Normal and Student t error distributions is able to recover the true state-dependent parameters for both low and high degree of dependence and all three considered SDs. Not surprisingly, the bias effect is quite small when we analyze the median levels (see column 3). As the quantile levels become more extreme (see columns 1, 2, 4 and 5), the ARB slightly increases but it still remains reasonably small. Also, under the ( $\mathcal{T}$ ) scenario the heavier tails of the Student t contribute to higher ARB and RMSE especially at the 10-th and 90-th percentiles.

To evaluate the classification performance of the proposed model, we report the average Adjusted Rand Index (ARI) of Hubert and Arabie 1985 and the misclassification rate (MCR). Specifically, we compare the classification obtained by the QHSMM with the one obtained from a QHMM under the assumption of a geometric SD. The state partition provided by the fitted models is obtained by taking the maximum,  $\max_{k \in \mathcal{S}} \gamma_{tk}$ , posteriori probability for every  $t = 1, \dots, T$ . The results in Table 3 show that when the true SD of the data generating process is geometric, the QHMM provides a slightly better classification both in terms of ARI and MCR (please compare the GEO row of Panel A with that of Panel C and the GEO row of Panel B with that of Panel D). This is not surprising as the QHMM implicitly assumes geometrically distributed sojourn distributions. The QHSMM, on the contrary, outperforms the QHMM in all other cases, as it can approximate arbitrarily well any SD and does not rely on a distributional assumption for  $d_k(u)$  (compare the SPO

**Table 3** Average and standard deviation (in brackets) values of the ARI and MCR for the QHSMM and QHMM under the three considered SDs and two distributions for the error term

| SD                     | $\tau$           |                  |                  |                  |                  |                  |                  |                  |                  |                  |
|------------------------|------------------|------------------|------------------|------------------|------------------|------------------|------------------|------------------|------------------|------------------|
|                        | (0.10, 0.10)     |                  | (0.25, 0.25)     |                  | (0.50, 0.50)     |                  | (0.75, 0.75)     |                  | (0.90, 0.90)     |                  |
|                        | ARI              | MCR              | ARI              | MCR              | ARI              | MCR              | ARI              | MCR              | ARI              | MCR              |
| <b>QHSMM</b>           |                  |                  |                  |                  |                  |                  |                  |                  |                  |                  |
| Panel A: $\mathcal{N}$ |                  |                  |                  |                  |                  |                  |                  |                  |                  |                  |
| SPO                    | 0.905<br>(0.022) | 0.024<br>(0.006) | 0.908<br>(0.021) | 0.023<br>(0.006) | 0.908<br>(0.021) | 0.023<br>(0.006) | 0.909<br>(0.021) | 0.023<br>(0.005) | 0.907<br>(0.022) | 0.024<br>(0.006) |
| SNB                    | 0.852<br>(0.026) | 0.037<br>(0.007) | 0.857<br>(0.026) | 0.036<br>(0.007) | 0.855<br>(0.026) | 0.036<br>(0.007) | 0.856<br>(0.026) | 0.036<br>(0.007) | 0.853<br>(0.026) | 0.037<br>(0.007) |
| GEO                    | 0.783<br>(0.028) | 0.057<br>(0.008) | 0.785<br>(0.029) | 0.057<br>(0.008) | 0.783<br>(0.028) | 0.057<br>(0.008) | 0.783<br>(0.029) | 0.057<br>(0.008) | 0.782<br>(0.028) | 0.058<br>(0.008) |
| Panel B: $\mathcal{T}$ |                  |                  |                  |                  |                  |                  |                  |                  |                  |                  |
| SPO                    | 0.874<br>(0.029) | 0.032<br>(0.008) | 0.885<br>(0.023) | 0.030<br>(0.006) | 0.885<br>(0.023) | 0.029<br>(0.006) | 0.885<br>(0.023) | 0.029<br>(0.006) | 0.876<br>(0.025) | 0.032<br>(0.007) |
| SNB                    | 0.810<br>(0.030) | 0.048<br>(0.008) | 0.824<br>(0.028) | 0.045<br>(0.008) | 0.824<br>(0.028) | 0.045<br>(0.008) | 0.824<br>(0.029) | 0.045<br>(0.008) | 0.811<br>(0.031) | 0.048<br>(0.008) |
| GEO                    | 0.737<br>(0.031) | 0.071<br>(0.009) | 0.745<br>(0.031) | 0.068<br>(0.009) | 0.744<br>(0.031) | 0.069<br>(0.009) | 0.743<br>(0.031) | 0.069<br>(0.009) | 0.735<br>(0.031) | 0.071<br>(0.009) |
| <b>QHMM</b>            |                  |                  |                  |                  |                  |                  |                  |                  |                  |                  |
| Panel C: $\mathcal{N}$ |                  |                  |                  |                  |                  |                  |                  |                  |                  |                  |
| SPO                    | 0.895<br>(0.022) | 0.027<br>(0.006) | 0.904<br>(0.021) | 0.024<br>(0.005) | 0.909<br>(0.020) | 0.023<br>(0.005) | 0.907<br>(0.020) | 0.024<br>(0.005) | 0.900<br>(0.021) | 0.025<br>(0.005) |
| SNB                    | 0.850<br>(0.026) | 0.038<br>(0.007) | 0.856<br>(0.025) | 0.036<br>(0.007) | 0.861<br>(0.025) | 0.035<br>(0.006) | 0.857<br>(0.025) | 0.036<br>(0.007) | 0.851<br>(0.026) | 0.037<br>(0.007) |
| GEO                    | 0.799<br>(0.026) | 0.053<br>(0.007) | 0.802<br>(0.026) | 0.052<br>(0.007) | 0.804<br>(0.026) | 0.052<br>(0.007) | 0.802<br>(0.027) | 0.052<br>(0.007) | 0.798<br>(0.027) | 0.053<br>(0.008) |
| Panel D: $\mathcal{T}$ |                  |                  |                  |                  |                  |                  |                  |                  |                  |                  |
| SPO                    | 0.856<br>(0.026) | 0.037<br>(0.007) | 0.874<br>(0.024) | 0.032<br>(0.006) | 0.883<br>(0.023) | 0.030<br>(0.006) | 0.876<br>(0.023) | 0.032<br>(0.006) | 0.858<br>(0.026) | 0.037<br>(0.007) |
| SNB                    | 0.803<br>(0.029) | 0.050<br>(0.008) | 0.820<br>(0.027) | 0.046<br>(0.007) | 0.830<br>(0.026) | 0.043<br>(0.007) | 0.822<br>(0.027) | 0.045<br>(0.007) | 0.806<br>(0.029) | 0.049<br>(0.008) |
| GEO                    | 0.752<br>(0.030) | 0.066<br>(0.009) | 0.762<br>(0.030) | 0.063<br>(0.008) | 0.768<br>(0.029) | 0.062<br>(0.008) | 0.763<br>(0.030) | 0.063<br>(0.008) | 0.752<br>(0.029) | 0.066<br>(0.008) |

and SNB rows of Panel A with those of Panel C, and the SPO and SNB rows of Panel B with those of Panel D), with very few exceptions at quantile levels  $\tau = (0.50, 0.50)$  and  $\tau = (0.75, 0.75)$  where the two models give comparable results.

We further evaluate the QHSMM introduced when a nonlinear quantile regression function of  $\mathbf{Y}$  given  $\mathbf{X}$  is considered. Similarly to Geraci (2019), the observations are generated from a two state HSMM using the following nonlinear quantile regression models. In the first scenario, we simulated the data from the following logistic model:

$$Y_t^{(j)} = \frac{\beta_{1,jk}}{1 + \exp((\beta_{2,jk} - X_t)/\beta_{3,jk})} + \epsilon_{tk}^{(j)}, \text{ for } j = 1, 2, \tag{31}$$

where the explanatory variable is drawn from a continuous uniform distribution,  $X_t \sim \mathcal{U}(0, 20)$ , and where  $Y_t^{(j)}$  and  $\epsilon_{tk}^{(j)}$  denote the  $j$ -th component of  $\mathbf{Y}_t$  and  $\epsilon_{tk}$ , respectively. For each component  $j$  and hidden state  $k$ , the true values of the parameters,  $\beta_{jk}$ , are given by  $\beta_{11} = (50, 12, 3)$ ,  $\beta_{21} = (10, 2, -1)$ ,  $\beta_{12} = (30, 5, 2)$  and  $\beta_{22} = (20, 11, -3)$ .

In the second scenario, following El Ghouch and Genton (2009) the data are generated according to the equation:

$$Y_t^{(j)} = \beta_{1,jk} + \beta_{2,jk} X_t^2 + \beta_{3,jk} X_t^3 + \gamma_j \exp(-4(X_t - 1)^2) + \epsilon_{tk}^{(j)}, \text{ for } j = 1, 2, \tag{32}$$

where  $X_t \sim \mathcal{U}(-1.1, 2.1)$  and where the parameter  $\gamma_j$  in (32) can be seen as a misspecification parameter that con-

trols the deviation from the polynomial function. As  $\gamma_j$  increases, the data structure becomes more complicated, and approximating the true curve by a polynomial becomes increasingly difficult. In this study we chose  $\gamma_1 = 4$  and  $\gamma_2 = 8$ . The true values of the parameters are given by  $\beta_{11} = (10, -6, 2.8)$ ,  $\beta_{21} = (3, 5, -0.5)$ ,  $\beta_{12} = (-3, 5, -2)$  and  $\beta_{22} = (8, 11, -3)$ .

For the error terms  $\epsilon_{tk}$  in (31) and (32), and the SDs we considered the same distributions adopted for the linear case. Examples of the simulated data are shown in Figures 1 and 2 from the two scenarios, respectively.

We fit the proposed QHSM for five quantile levels, i.e.,  $\tau = (0.10, 0.10)$ ,  $\tau = (0.25, 0.25)$ ,  $\tau = (0.50, 0.50)$ ,  $\tau = (0.75, 0.75)$  and  $\tau = (0.90, 0.90)$ . For each model, we report the Proportion of Negative Residuals (PNR):

$$PNR(\tau_j) = \frac{1}{T} \sum_{t=1}^T \mathbf{1}(Y_t^{(j)} < \widehat{Q}_{Y_t^{(j)}|X_t}), \text{ for } j = 1, 2, \tag{33}$$

with  $\widehat{Q}_{Y_t^{(j)}|X_t}$  being the fitted conditional quantile of  $Y_t^{(j)}$  at level  $\tau_j$ ,  $j = 1, 2$ . If the model is correctly specified, the PNR should be approximately equal to  $\tau_j$  for each outcome. The results contained in Tables 4 and 5 are averaged over  $B = 1000$  Monte Carlo replications.

By looking at the results, PNR rates are in general coherent with the selected quantile level. Further, one can observe that PNRs are typically closer to the nominal values  $\tau$  in the first scenario as opposed to the second one. This may be explained by the fact that the fitted model provides a better linear approximation of the true logistic quantile regression function than the more complex nonlinear model in the second scenario. Moreover, in this latter case, the PNRs obtained for the second component  $Y_2$  do worse than the corresponding ones for  $Y_1$ , which is mainly due to the relatively larger variance associated to the misspecification parameter  $\gamma_j$  in (32). Overall, in both cases the PNRs are slightly, especially at the tails, below (at  $\tau = (0.10, 0.10)$ ) or above (at  $\tau = (0.90, 0.90)$ ) the expected proportions.

Finally, to assess the performance of penalized likelihood criteria (AIC, BIC and ICL) for selecting the number of hidden states and maximum sojourn times, we considered the same simulation experiment where the sojourns are generated from a beta-binomial distribution:

$$d_k(u) = \binom{U_k}{u} \frac{\mathcal{B}(u + \alpha_k, U - u + \lambda_k)}{\mathcal{B}(\alpha_k, \lambda_k)}, \text{ } u = 1, 2, \dots, U_k, \tag{34}$$

where  $\mathcal{B}(\cdot, \cdot)$  is the beta function,  $\alpha_1 = 6$ ,  $\alpha_2 = 0.7$ ,  $\lambda_1 = 3$ ,  $\lambda_2 = 2$ ,  $U_1 = 20$  and  $U_2 = 10$ . For each of the simulated  $B = 100$  datasets, we fit the QHSM with  $K = 2, 3, 4$

over the grid of state-dependent maximum sojourn times  $(10, 15, 20) \times \dots \times (10, 15, 20)$ , and select the best combination of  $(K, U)$  associated to the lowest penalized likelihood criteria. Table 6 reports the number of times each criterion correctly identifies the number of latent states and maximum support points of the SDs (Panel A) and the absolute frequency distributions of the selected  $K$  for each of the three criteria (Panel B).

As one can see in Panel A, all three criteria work relatively well at  $\tau = (0.50, 0.50)$ . By contrast, as we move towards the tails of the distribution of the responses, the ICL outperforms both the AIC and BIC and correctly identifies the pair  $(K, U)$  that was used to generate the data more than 96% across all simulation scenarios. Moving onto Panel B, the AIC consistently performs worse than the BIC but both mostly overestimate the true number of states. These results suggest that regardless of the distribution on the error terms and SDs, the ICL yields superior performance and captures serial heterogeneity in the data in a more parsimonious manner compared to the other criteria, easing the interpretation of the latent states.

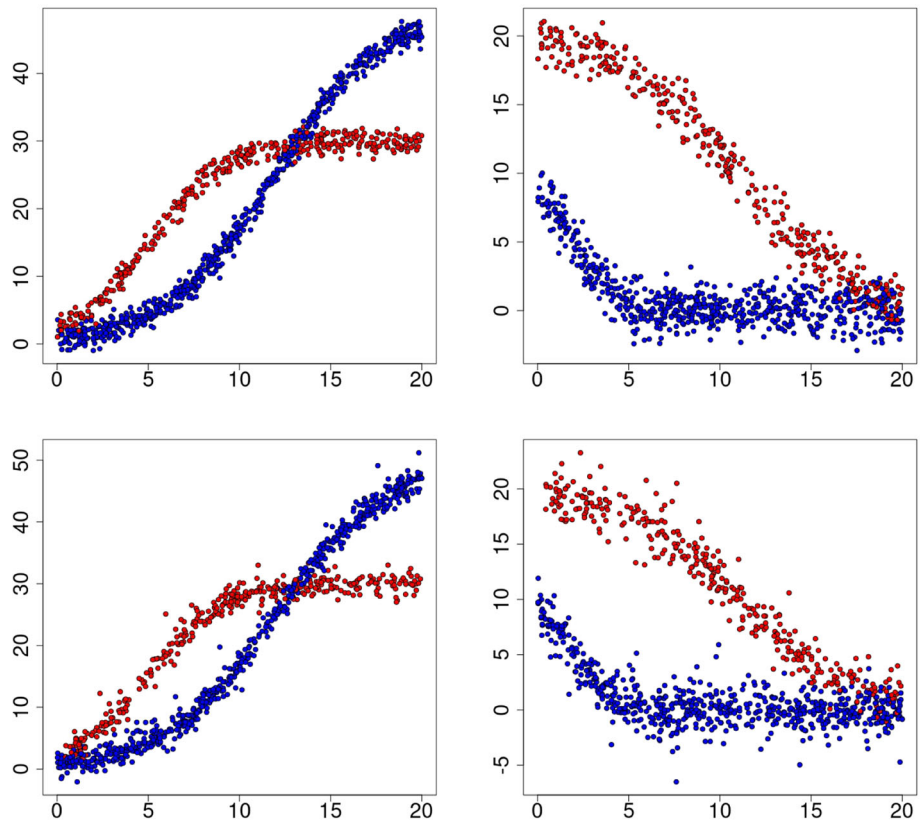
## 5 Application

In this section we apply the proposed methodology to air pollution data collected by the Lazio Regional Agency for Environmental Prevention and Protection (ARPA Lazio, <https://www.arpalazio.it>) in Italy. The ARPA Lazio provides information regarding the regional state of the environment and environmental trends, performing scientific, technical and research functions as well as assessment, monitoring, control and supporting local and health authorities. The time series used in this research are freely available from the ARPA Lazio website (<https://www.arpalazio.net/main/aria/sci/basedati/chimici/chimici.php>).

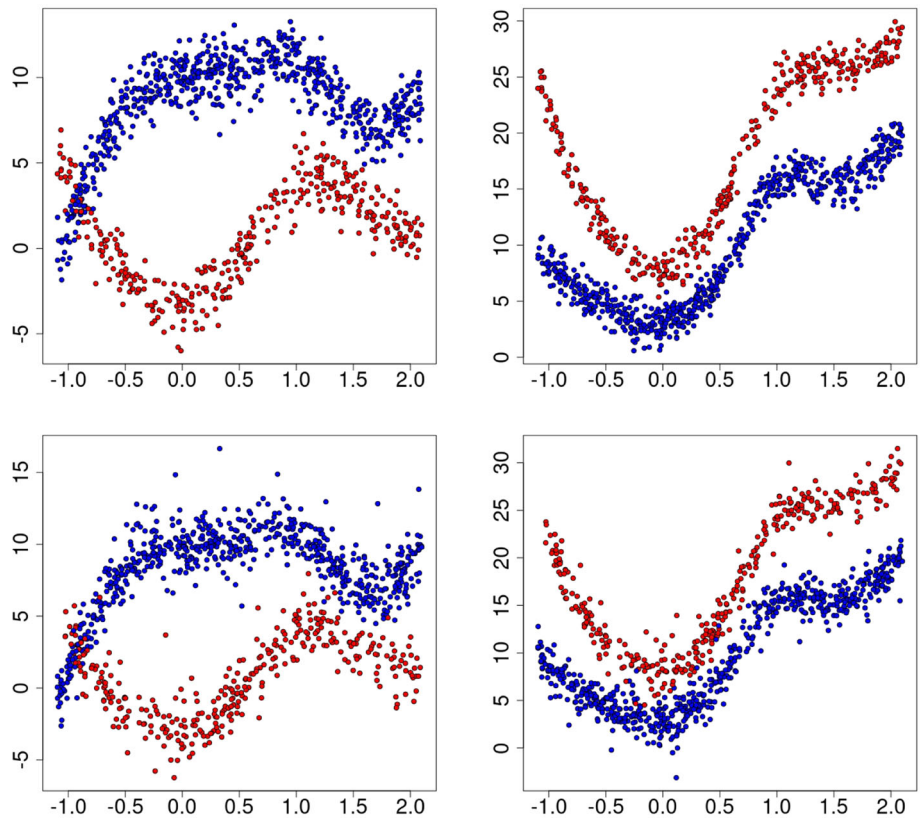
### 5.1 Data description

The data considered originate from a regional monitoring network system developed by the Lazio Region and have already been discussed in the work of Maruotti et al. (2017). This system has been organized in order to respond to the increasing demand for environmental information, but also for providing reliable data suitable for policy-related aspects. The network recorded concentrations of nine air pollutants on hourly basis at monitoring stations in the central area of the city of Rieti, Italy. The Rieti site has been chosen as it is classified as a traffic location, although it is located at a short distance from green areas and forests that facilitate the movement of air masses and removal of pollutants. In this work we consider three major air pollutants, i.e., Particulate Matter with aerodynamic diameter less than  $2.5 \mu\text{m}$  ( $\text{PM}_{2.5}$ ),

**Fig. 1** Examples of data generated from the first scenario. Scatterplot of  $Y_1$  (first column) and  $Y_2$  (second column) under normal (first row) and Student t (second row) errors as a function of the included covariate, when using shifted-Poisson SDs. Red and blue data points distinguish the two latent states



**Fig. 2** Examples of data generated from the second scenario. Scatterplot of  $Y_1$  (first column) and  $Y_2$  (second column) under normal (first row) and Student t (second row) errors as a function of the included covariate, when using shifted-Poisson SDs. Red and blue data points distinguish the two latent states



**Table 4** Average and standard deviation (in brackets) values of the PNR for  $Y_1$  and  $Y_2$  under the considered SDs and error term distributions for the first scenario over 1000 samples

| PNR               | $\tau$           |                  |                  |                  |                  |                  |                  |                  |                  |                  |
|-------------------|------------------|------------------|------------------|------------------|------------------|------------------|------------------|------------------|------------------|------------------|
|                   | (0.10, 0.10)     |                  | (0.25, 0.25)     |                  | (0.50, 0.50)     |                  | (0.75, 0.75)     |                  | (0.90, 0.90)     |                  |
| <i>Panel A: N</i> |                  |                  |                  |                  |                  |                  |                  |                  |                  |                  |
| SPO               | 0.096<br>(0.005) | 0.083<br>(0.004) | 0.265<br>(0.008) | 0.233<br>(0.007) | 0.497<br>(0.007) | 0.504<br>(0.008) | 0.751<br>(0.006) | 0.769<br>(0.006) | 0.920<br>(0.005) | 0.912<br>(0.004) |
| SNB               | 0.091<br>(0.007) | 0.079<br>(0.006) | 0.251<br>(0.010) | 0.224<br>(0.009) | 0.492<br>(0.012) | 0.501<br>(0.014) | 0.750<br>(0.007) | 0.770<br>(0.006) | 0.918<br>(0.005) | 0.911<br>(0.004) |
| GEO               | 0.088<br>(0.008) | 0.076<br>(0.006) | 0.246<br>(0.012) | 0.231<br>(0.017) | 0.499<br>(0.007) | 0.502<br>(0.008) | 0.753<br>(0.007) | 0.769<br>(0.006) | 0.920<br>(0.004) | 0.913<br>(0.005) |
| <i>Panel B: T</i> |                  |                  |                  |                  |                  |                  |                  |                  |                  |                  |
| SPO               | 0.097<br>(0.006) | 0.080<br>(0.004) | 0.266<br>(0.007) | 0.232<br>(0.007) | 0.498<br>(0.007) | 0.502<br>(0.008) | 0.750<br>(0.007) | 0.768<br>(0.006) | 0.918<br>(0.005) | 0.914<br>(0.004) |
| SNB               | 0.089<br>(0.006) | 0.080<br>(0.005) | 0.256<br>(0.011) | 0.224<br>(0.010) | 0.495<br>(0.010) | 0.502<br>(0.012) | 0.750<br>(0.007) | 0.768<br>(0.006) | 0.917<br>(0.005) | 0.912<br>(0.004) |
| GEO               | 0.095<br>(0.008) | 0.080<br>(0.005) | 0.249<br>(0.014) | 0.231<br>(0.017) | 0.500<br>(0.007) | 0.500<br>(0.009) | 0.752<br>(0.007) | 0.769<br>(0.006) | 0.918<br>(0.004) | 0.914<br>(0.005) |

**Table 5** Average and standard deviation (in brackets) values of the PNR for  $Y_1$  and  $Y_2$  under the considered SDs and error term distributions for the second scenario over 1000 samples

| PNR               | $\tau$           |                  |                  |                  |                  |                  |                  |                  |                  |                  |
|-------------------|------------------|------------------|------------------|------------------|------------------|------------------|------------------|------------------|------------------|------------------|
|                   | (0.10, 0.10)     |                  | (0.25, 0.25)     |                  | (0.50, 0.50)     |                  | (0.75, 0.75)     |                  | (0.90, 0.90)     |                  |
| <i>Panel A: N</i> |                  |                  |                  |                  |                  |                  |                  |                  |                  |                  |
| SPO               | 0.095<br>(0.010) | 0.090<br>(0.017) | 0.249<br>(0.007) | 0.239<br>(0.007) | 0.479<br>(0.008) | 0.482<br>(0.008) | 0.743<br>(0.008) | 0.760<br>(0.008) | 0.907<br>(0.033) | 0.923<br>(0.011) |
| SNB               | 0.090<br>(0.008) | 0.083<br>(0.013) | 0.250<br>(0.006) | 0.239<br>(0.006) | 0.479<br>(0.008) | 0.481<br>(0.008) | 0.742<br>(0.008) | 0.758<br>(0.008) | 0.922<br>(0.024) | 0.926<br>(0.009) |
| GEO               | 0.092<br>(0.008) | 0.084<br>(0.014) | 0.251<br>(0.006) | 0.238<br>(0.006) | 0.486<br>(0.008) | 0.483<br>(0.009) | 0.742<br>(0.008) | 0.758<br>(0.008) | 0.909<br>(0.042) | 0.926<br>(0.011) |
| <i>Panel B: T</i> |                  |                  |                  |                  |                  |                  |                  |                  |                  |                  |
| SPO               | 0.090<br>(0.009) | 0.082<br>(0.015) | 0.254<br>(0.006) | 0.245<br>(0.007) | 0.483<br>(0.008) | 0.485<br>(0.008) | 0.746<br>(0.008) | 0.756<br>(0.008) | 0.933<br>(0.012) | 0.926<br>(0.009) |
| SNB               | 0.088<br>(0.008) | 0.079<br>(0.012) | 0.254<br>(0.006) | 0.244<br>(0.007) | 0.482<br>(0.008) | 0.484<br>(0.008) | 0.745<br>(0.008) | 0.754<br>(0.007) | 0.935<br>(0.011) | 0.925<br>(0.009) |
| GEO               | 0.090<br>(0.008) | 0.080<br>(0.012) | 0.254<br>(0.006) | 0.243<br>(0.007) | 0.487<br>(0.008) | 0.485<br>(0.008) | 0.745<br>(0.008) | 0.755<br>(0.007) | 0.923<br>(0.031) | 0.924<br>(0.011) |

Ozone ( $O_3$ ) and Nitrogen Dioxide ( $NO_2$ ) from January 01, 2019 to June 14, 2021. We averaged the pollution data to daily frequency and all concentrations are expressed in  $\mu g/m^3$ .

Atmospheric variables also play a major role in determining the level of exposure to particular pollutants and capturing time dependence as proxies for seasonal variations and characteristics. Since pollution episodes are triggered by specific atmospheric factors, we include the following variables, namely the daily average wind speed, temperature, pressure and humidity. Table 7 presents the main descriptive

statistics for the response variables and the set of included predictors. The asymmetry in the distributions of the pollutants is noted by examining the mean-median relationship and the five summary statistics indicate severe departures from the Gaussian assumptions, presenting high kurtosis and outlying values. In the same table we also report the empirical correlation coefficients between each response variable which clearly highlight a positive correlation between  $PM_{25}$  and  $NO_2$ , and a negative association with  $O_3$ . Therefore, the dependence structure among different pollutants, which can-

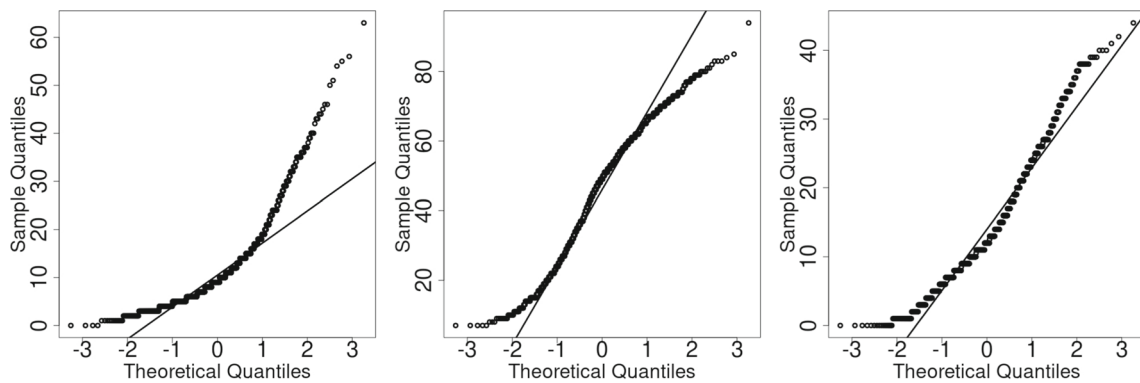


**Table 6** Number of correctly identified hidden states  $K$  and maximum sojourn times  $U$  (Panel A) and absolute frequency distribution of the selected number of states (Panel B) under Normal and Student  $t$  errors over 100 replications

|                          | $\tau$       |     |     |              |     |     |              |     |     |              |     |     |              |     |     |
|--------------------------|--------------|-----|-----|--------------|-----|-----|--------------|-----|-----|--------------|-----|-----|--------------|-----|-----|
|                          | (0.10, 0.10) |     |     | (0.25, 0.25) |     |     | (0.50, 0.50) |     |     | (0.75, 0.75) |     |     | (0.90, 0.90) |     |     |
|                          | AIC          | BIC | ICL | AIC          | BIC | ICL | AIC          | BIC | ICL | AIC          | BIC | ICL | AIC          | BIC | ICL |
| Panel A                  |              |     |     |              |     |     |              |     |     |              |     |     |              |     |     |
| $\mathcal{N}$            | 0            | 1   | 100 | 0            | 54  | 99  | 30           | 98  | 99  | 0            | 61  | 100 | 0            | 2   | 99  |
| $\mathcal{T}$            | 0            | 0   | 82  | 0            | 27  | 97  | 72           | 86  | 97  | 0            | 22  | 97  | 0            | 0   | 92  |
| Panel B                  |              |     |     |              |     |     |              |     |     |              |     |     |              |     |     |
| Panel B.1: $\mathcal{N}$ |              |     |     |              |     |     |              |     |     |              |     |     |              |     |     |
| 2                        | 0            | 1   | 100 | 0            | 56  | 100 | 30           | 99  | 99  | 0            | 62  | 100 | 0            | 2   | 100 |
| 3                        | 0            | 70  | 0   | 0            | 44  | 0   | 4            | 1   | 1   | 0            | 38  | 0   | 0            | 80  | 0   |
| 4                        | 100          | 29  | 0   | 100          | 0   | 0   | 66           | 0   | 0   | 100          | 0   | 0   | 100          | 18  | 0   |
| Panel B.2: $\mathcal{T}$ |              |     |     |              |     |     |              |     |     |              |     |     |              |     |     |
| 2                        | 0            | 0   | 88  | 0            | 32  | 100 | 72           | 99  | 100 | 0            | 28  | 100 | 0            | 0   | 98  |
| 3                        | 0            | 15  | 11  | 1            | 68  | 0   | 13           | 1   | 0   | 0            | 72  | 0   | 0            | 56  | 2   |
| 4                        | 100          | 85  | 1   | 99           | 0   | 0   | 15           | 0   | 0   | 100          | 0   | 0   | 100          | 44  | 0   |

**Table 7** Summary statistics of the sample data

| Variable           | Minimum          | 1st quartile   | Median          | Mean     | 3rd quartile | Maximum |
|--------------------|------------------|----------------|-----------------|----------|--------------|---------|
| PM <sub>25</sub>   | 0                | 6              | 9               | 11.846   | 15           | 63      |
| O <sub>3</sub>     | 7                | 31             | 49              | 46.340   | 61           | 94      |
| NO <sub>2</sub>    | 0                | 8              | 12              | 14.139   | 20           | 44      |
| Wind Speed         | 0                | 6              | 8               | 8.839    | 11           | 29      |
| Temperature        | 0                | 11             | 16              | 16.659   | 22           | 34      |
| Pressure           | 989              | 1011           | 1015            | 1015.127 | 1020         | 1036    |
| Humidity           | 17               | 42             | 56              | 58.415   | 74           | 100     |
| Correlation matrix |                  |                |                 |          |              |         |
|                    | PM <sub>25</sub> | O <sub>3</sub> | NO <sub>2</sub> |          |              |         |
| PM <sub>25</sub>   | 1                |                |                 |          |              |         |
| O <sub>3</sub>     | -0.566           | 1              |                 |          |              |         |
| NO <sub>2</sub>    | 0.743            | -0.656         | 1               |          |              |         |



**Fig. 3** From left to right, univariate normal QQ plots for PM<sub>25</sub>, O<sub>3</sub> and NO<sub>2</sub>

**Table 8** Log-likelihood, AIC, BIC and ICL values for a varying number of hidden states. Bold font highlights the best values for the considered criteria (lower-is-better)

| K | $\tau$   |                 |                 |                 | (0.50, 0.50, 0.50) |                 |                 |                 | (0.75, 0.75, 0.75) |                 |          |          | (0.90, 0.90, 0.90) |          |                 |                 |
|---|----------|-----------------|-----------------|-----------------|--------------------|-----------------|-----------------|-----------------|--------------------|-----------------|----------|----------|--------------------|----------|-----------------|-----------------|
|   | Loglik   | AIC             | BIC             | ICL             | Loglik             | AIC             | BIC             | ICL             | Loglik             | AIC             | BIC      | ICL      | Loglik             | AIC      | BIC             | ICL             |
| 1 | -8284.99 | 16611.98        | 16712.73        | 16712.73        | -8553.92           | 17149.85        | 17250.61        | 17250.61        | -9002.05           | 18046.10        | 18146.86 | 18146.86 | -8415.43           | 17040.86 | 17544.65        | 17624.07        |
| 2 | -7893.43 | 15986.87        | 16466.66        | 16521.09        | -8093.62           | 16387.25        | 16867.04        | 16937.16        | -8149.87           | 16635.73        | 17441.79 | 17490.93 | -7926.16           | 16258.32 | <b>17232.30</b> | <b>17396.04</b> |
| 3 | -7656.82 | 15549.64        | <b>16115.80</b> | <b>16132.46</b> | -7808.18           | 15952.36        | 16758.42        | 16801.01        | -7820.87           | <b>16181.74</b> | 17477.18 | 17638.80 | -7761.82           | 16191.63 | 17794.14        | 17916.22        |
| 4 | -7434.64 | 15315.28        | 16385.22        | 16424.51        | -7655.14           | 15756.27        | 16826.21        | 16908.00        | -7400.95           | 15509.90        | 17208.37 | 17284.46 | -7400.95           | 15509.90 | 17208.37        | 17284.46        |
| 5 | -7311.89 | <b>15083.79</b> | 16187.32        | 16240.48        | -7523.67           | <b>15507.34</b> | <b>16610.86</b> | <b>16755.08</b> | -7400.95           | 15509.90        | 17208.37 | 17284.46 | -7400.95           | 15509.90 | 17208.37        | 17284.46        |
| 6 | -7195.51 | 15099.02        | 16797.49        | 16854.63        | -7400.95           | 15509.90        | 17208.37        | 17284.46        | -7400.95           | 15509.90        | 17208.37 | 17284.46 | -7400.95           | 15509.90 | 17208.37        | 17284.46        |

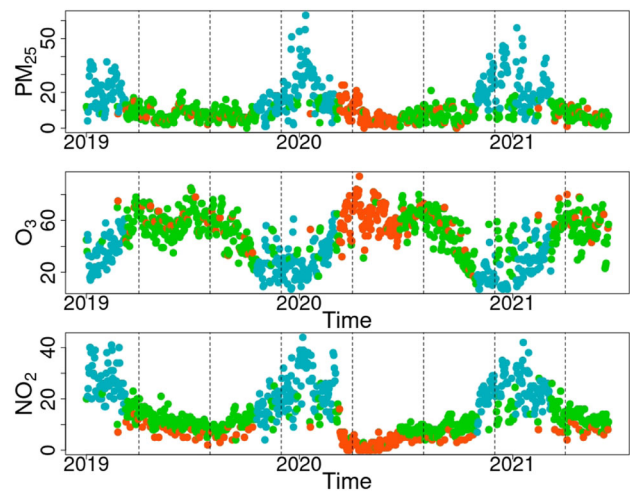
not be detected by univariate methods, constitutes a crucial aspect of the analysis and should not be neglected.

From a graphical standpoint, Fig. 3 shows the normal QQ plots for the PM<sub>25</sub>, O<sub>3</sub> and NO<sub>2</sub> time series. These reveal the presence of potentially influential observations in the data, heavy tails and skewness for all three outcomes. This exploratory analysis and preliminary considerations motivate us to consider a joint quantile regression approach as investigative tool.

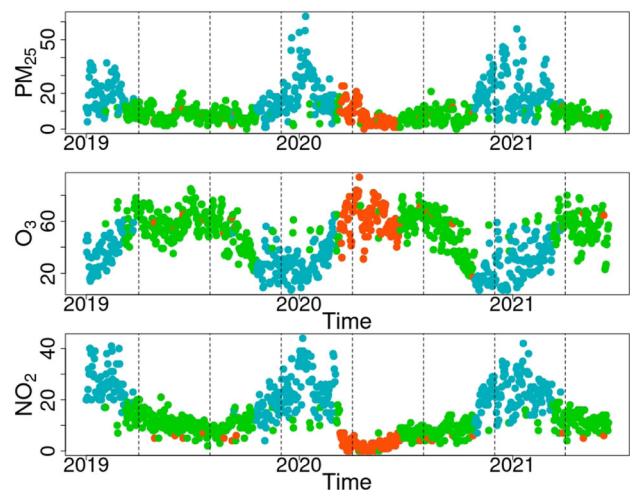
## 5.2 Results

We jointly model the concentrations of PM<sub>25</sub>, O<sub>3</sub> and NO<sub>2</sub> as a function of Wind Speed, Temperature, Pressure and Humidity at quantile levels  $\tau = (0.50, 0.50, 0.50)$ ,  $\tau = (0.75, 0.75, 0.75)$  and  $\tau = (0.90, 0.90, 0.90)$ . Considering the 75-th and 90-th percentiles puts emphasis on alert thresholds for ambient air pollution associated with high levels of concentrations of chemicals. As a first step of the analysis, we fit the proposed QHSMM for a sequence of states  $K$  from 1 to 6 and a  $K$ -dimensional grid,  $\mathcal{U} \subset \mathbb{R}_{\geq 0}^K$ , of maximum sojourn times,  $(10, 15, \dots, 60) \times \dots \times (10, 15, \dots, 60)$ . To select the optimal combination of  $K$  and  $\mathbf{U}$ , with  $\mathbf{U} \in \mathcal{U}$ , and avoid the computational cost of a full grid search over  $\mathcal{U}$ , we employ the greedy search algorithm described in Sect. 3.2 with 200 starting points. Table 8 reports the log-likelihood, AIC, BIC and ICL values for the fitted models at the investigated quantile levels. The AIC selects  $K = 5$  states for all quantile levels meanwhile the BIC and ICL are more aligned in the choice of  $K$  as they identify 3, 5 and 4 states for  $\tau = (0.50, 0.50, 0.50)$ ,  $\tau = (0.75, 0.75, 0.75)$  and  $\tau = (0.90, 0.90, 0.90)$ . This is not surprising since the AIC tends to overestimate the number of hidden states and, for this reason, we will not consider it hereafter. We can also see that the ICL values associated to  $K = 3$  and  $K = 5$  are extremely similar at  $\tau = (0.75, 0.75, 0.75)$ . Following these considerations and looking at Table 8, we select the best fitted models with  $K$  equal to 3, 3 and 4 according to both the BIC and ICL criteria at  $\tau = (0.50, 0.50, 0.50)$ ,  $\tau = (0.75, 0.75, 0.75)$  and  $\tau = (0.90, 0.90, 0.90)$ , respectively.

Figures 4, 5 and 6 report the classification results according to the selected models at the investigated quantile levels. Each plot shows the data points colored according to the estimated posterior probability of class membership,  $\hat{\gamma}_{tk}$ , with the vertical lines separating blocks of four months. In our study, the latent components can be associated to specific exposure regimes characterized by seasonal weather conditions. Specifically, blue points (state 2) tend to cluster days in late autumn, winter, and early spring, while green ones (state 3) are generally inferred during late spring and summer. At  $\tau = (0.90, 0.90, 0.90)$  (see Figure 6), the four state QHSMM identifies a similar classification pattern with violet points (state 4) occurring mainly from spring until the end of

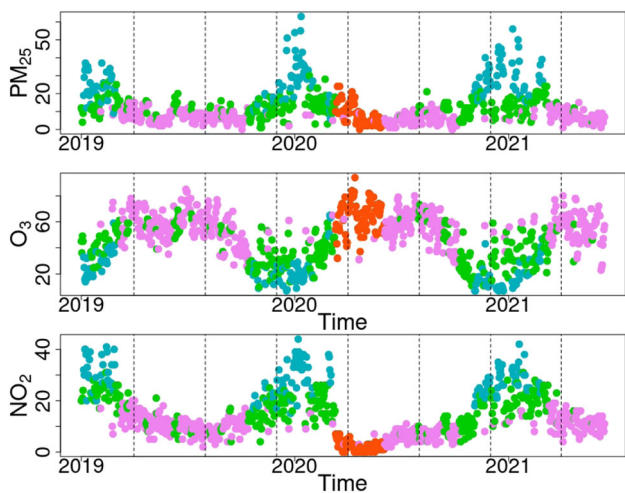


**Fig. 4** Time series classified according to their estimated posterior probability of class membership at  $\tau = (0.50, 0.50, 0.50)$ . Vertical lines separate blocks of four months



**Fig. 5** Time series classified according to their estimated posterior probability of class membership at  $\tau = (0.75, 0.75, 0.75)$ . Vertical lines separate blocks of four months

summer, but also sporadically during the rest of the year. The overall results reflect the seasonal variation in atmospheric pollutants, with high concentrations of PM<sub>25</sub> and NO<sub>2</sub> in winter and high O<sub>3</sub> concentrations in summer and warm months. It is also worth noting that the minimum values of particulate matter and nitrogen dioxide, and the maximum level of ozone were reached in state 1 (orange dots) during the period March–June 2020. These sudden changes may be possibly due to the implementation of lockdown measures to contain the COVID-19 outbreak in Italy (Bassani et al. 2021; Putaud et al. 2021). Indeed, in the first weeks of March, significant PM<sub>25</sub> and NO<sub>2</sub> declines were observed which led to O<sub>3</sub> peaks resulting from reduced titration with nitrogen oxides. As restrictions were lifted in May–June, chemical concentrations settled again around the 2019 and early 2021 levels.



**Fig. 6** Time series classified according to their estimated posterior probability of class membership at  $\tau = (0.90, 0.90, 0.90)$ . Vertical lines separate blocks of four months

The estimated transition probability matrices of the latent semi-Markov chain (see Table 9) confirm that pollutant concentrations alternate between good and poor air conditions. The off-diagonal elements demonstrate that the hidden process stays in state 3 (low levels of particulate matter, nitrogen dioxide and moderate levels of ozone), and 4 in the case  $\tau = (0.90, 0.90, 0.90)$ , with temporary changes towards more severe PM<sub>25</sub> and NO<sub>2</sub> concentrations (state 2) or higher ozone episodes (state 1). Meanwhile, direct transitions between states 1 and 2 are very unlikely at all three quantile levels.

Taking the effect of covariates into account, Table 10 shows the estimated state-specific regression parameters,  $\beta_k, k = 1, \dots, K$ , at each quantile level, respectively. Standard errors are computed via parametric bootstrap using  $H = 1000$  resamples as illustrated in Section 3 and point

estimates are displayed in boldface when significant at the standard 5% level.

The estimated effects of the included covariates tend to be nonlinear, state-specific and are generally more pronounced in the upper end of the distribution of the responses. States 1 and 3 yield similar estimates as they are both associated with the best air conditions, especially when looking at effect of pressure and humidity. This is in sharp contrast with the point estimates in state 2 which is characterized by hazardous air quality. Among the four factors, wind speed and temperature exert the strongest influence on the considered pollutants in all seasons. In particular, PM<sub>25</sub> and NO<sub>2</sub> are negatively associated with both variables and such association increases during the coldest months of the year. Wind intensity, therefore, contributes considerably to the reduction of pollution, in particular in at-risk situations as identified by state 2. On the other hand, O<sub>3</sub> concentration is positively associated with wind speed meanwhile the effect of temperature is positive in late-spring and hot weather, but negative throughout the rest of the year. Humidity can also help to decrease ozone pollution because the moisture in the air could enhance the condensation of water and slow down ozone production, but also reduce PM<sub>25</sub> and NO<sub>2</sub>. Further, the concentrations of PM<sub>25</sub> and NO<sub>2</sub> are positively associated with atmospheric pressure and negatively associated with O<sub>3</sub>.

We conclude the analysis by reporting the estimated correlation matrices,  $\Psi_k, k = 1, \dots, K$ , (see Table 11) in order to provide an indirect measure of tail dependence between the outcomes. Firstly, regardless of the state, the correlation coefficients between the air pollutants are generally significant, indicating that in this case fitting univariate quantile regressions separately would be inappropriate. Secondly, the correlation coefficients are increasing somewhat with  $\tau$ . Finally, the estimates depict a data correlation structure that varies among the latent groups as the correlation between

**Table 9** Estimated transition probabilities for different quantile levels

| states                      | 1             | 2             | 3             | 4             |
|-----------------------------|---------------|---------------|---------------|---------------|
| Panel A: (0.50, 0.50, 0.50) |               |               |               |               |
| 1                           | 0             | 0.028 (0.009) | 0.972 (0.009) |               |
| 2                           | 0.000 (0.001) | 0             | 1.000 (0.001) |               |
| 3                           | 0.836 (0.020) | 0.164 (0.020) | 0             |               |
| Panel B: (0.75, 0.75, 0.75) |               |               |               |               |
| 1                           | 0             | 0.000 (0.007) | 1.000 (0.007) |               |
| 2                           | 0.000 (0.000) | 0             | 1.000 (0.000) |               |
| 3                           | 0.543 (0.030) | 0.457 (0.030) | 0             |               |
| Panel C: (0.90, 0.90, 0.90) |               |               |               |               |
| 1                           | 0             | 0.000 (0.000) | 0.000 (0.009) | 1.000 (0.009) |
| 2                           | 0.000 (0.000) | 0             | 0.966 (0.014) | 0.034 (0.014) |
| 3                           | 0.016 (0.007) | 0.482 (0.029) | 0             | 0.502 (0.030) |
| 4                           | 0.088 (0.025) | 0.097 (0.021) | 0.815 (0.032) | 0             |

**Table 10** Regression parameters estimates. Point estimates are displayed in boldface when significant at the standard 5% level

| States                      | 1                           |                            |                             | 2                            |                             |                             | 3                           |                            |                            | 4                           |                            |                            |
|-----------------------------|-----------------------------|----------------------------|-----------------------------|------------------------------|-----------------------------|-----------------------------|-----------------------------|----------------------------|----------------------------|-----------------------------|----------------------------|----------------------------|
|                             | PM <sub>2.5</sub>           | O <sub>3</sub>             | NO <sub>2</sub>             | PM <sub>2.5</sub>            | O <sub>3</sub>              | NO <sub>2</sub>             | PM <sub>2.5</sub>           | O <sub>3</sub>             | NO <sub>2</sub>            | PM <sub>2.5</sub>           | O <sub>3</sub>             | NO <sub>2</sub>            |
| Panel A: (0.50, 0.50, 0.50) |                             |                            |                             |                              |                             |                             |                             |                            |                            |                             |                            |                            |
| Intercept                   | <b>-192.006</b><br>(28.810) | <b>408.316</b><br>(58.965) | <b>-166.952</b><br>(20.840) | <b>-488.373</b><br>(130.720) | 149.664<br>(157.563)        | -110.794<br>(101.024)       | <b>-163.261</b><br>(20.776) | <b>690.359</b><br>(55.920) | <b>-52.286</b><br>(17.086) | <b>-148.541</b><br>(22.284) | <b>575.923</b><br>(59.029) | <b>-52.309</b><br>(19.866) |
| Wind Speed                  | <b>-0.323</b><br>(0.045)    | <b>0.320</b><br>(0.093)    | <b>-0.112</b><br>(0.034)    | <b>-0.538</b><br>(0.209)     | <b>1.869</b><br>(0.256)     | -0.227<br>(0.160)           | -0.059<br>(0.035)           | <b>0.704</b><br>(0.086)    | <b>-0.196</b><br>(0.028)   | -0.055<br>(0.034)           | <b>0.492</b><br>(0.091)    | <b>-0.278</b><br>(0.031)   |
| Temperature                 | <b>-0.236</b><br>(0.025)    | <b>-0.346</b><br>(0.050)   | <b>0.042</b><br>(0.019)     | <b>-1.095</b><br>(0.117)     | 0.237<br>(0.145)            | <b>-1.075</b><br>(0.089)    | <b>-0.100</b><br>(0.020)    | <b>0.314</b><br>(0.048)    | <b>-0.425</b><br>(0.015)   | <b>-0.149</b><br>(0.019)    | <b>-0.176</b><br>(0.049)   | <b>-0.428</b><br>(0.017)   |
| Pressure                    | <b>0.209</b><br>(0.028)     | <b>-0.309</b><br>(0.057)   | <b>0.167</b><br>(0.020)     | <b>0.515</b><br>(0.126)      | -0.118<br>(0.152)           | 0.156<br>(0.097)            | <b>0.173</b><br>(0.020)     | <b>-0.610</b><br>(0.054)   | <b>0.076</b><br>(0.016)    | 0.009<br>(0.022)            | <b>0.165</b><br>(0.021)    | <b>0.081</b><br>(0.019)    |
| Humidity                    | <b>-0.105</b><br>(0.012)    | <b>-0.619</b><br>(0.023)   | <b>0.042</b><br>(0.008)     | -0.006<br>(0.053)            | <b>-0.248</b><br>(0.066)    | <b>-0.155</b><br>(0.041)    | <b>-0.030</b><br>(0.009)    | <b>-0.603</b><br>(0.022)   | <b>-0.059</b><br>(0.007)   | <b>-0.063</b><br>(0.009)    | <b>-0.582</b><br>(0.024)   | <b>-0.069</b><br>(0.008)   |
| Panel B: (0.75, 0.75, 0.75) |                             |                            |                             |                              |                             |                             |                             |                            |                            |                             |                            |                            |
| Intercept                   | <b>-100.301</b><br>(37.049) | <b>151.191</b><br>(75.844) | <b>-119.282</b><br>(17.254) | <b>-482.680</b><br>(75.045)  | <b>174.791</b><br>(80.288)  | <b>-108.564</b><br>(53.198) | <b>-176.381</b><br>(21.221) | <b>778.126</b><br>(50.143) | <b>-44.040</b><br>(15.392) | <b>-148.541</b><br>(22.284) | <b>575.923</b><br>(59.029) | <b>-52.309</b><br>(19.866) |
| Wind Speed                  | <b>-0.816</b><br>(0.059)    | -0.022<br>(0.115)          | <b>-0.104</b><br>(0.027)    | <b>-0.584</b><br>(0.124)     | <b>1.832</b><br>(0.126)     | <b>-0.270</b><br>(0.085)    | <b>-0.069</b><br>(0.033)    | <b>0.554</b><br>(0.079)    | <b>-0.229</b><br>(0.024)   | -0.055<br>(0.034)           | <b>0.492</b><br>(0.091)    | <b>-0.278</b><br>(0.031)   |
| Temperature                 | <b>-0.621</b><br>(0.032)    | <b>-0.723</b><br>(0.062)   | <b>0.083</b><br>(0.015)     | <b>-1.170</b><br>(0.067)     | <b>0.212</b><br>(0.071)     | <b>-1.118</b><br>(0.046)    | <b>-0.105</b><br>(0.018)    | <b>-0.153</b><br>(0.045)   | <b>-0.412</b><br>(0.014)   | <b>-0.149</b><br>(0.019)    | <b>-0.176</b><br>(0.049)   | <b>-0.428</b><br>(0.017)   |
| Pressure                    | <b>0.137</b><br>(0.036)     | -0.039<br>(0.073)          | <b>0.119</b><br>(0.017)     | <b>0.518</b><br>(0.072)      | -0.132<br>(0.077)           | <b>0.158</b><br>(0.051)     | <b>0.190</b><br>(0.020)     | <b>-0.675</b><br>(0.048)   | <b>0.070</b><br>(0.015)    | 0.009<br>(0.022)            | <b>0.165</b><br>(0.021)    | <b>0.081</b><br>(0.019)    |
| Humidity                    | <b>-0.151</b><br>(0.014)    | <b>-0.599</b><br>(0.029)   | <b>0.058</b><br>(0.007)     | -0.053<br>(0.029)            | <b>-0.294</b><br>(0.032)    | <b>-0.170</b><br>(0.021)    | <b>-0.053</b><br>(0.009)    | <b>-0.646</b><br>(0.021)   | <b>-0.061</b><br>(0.007)   | <b>-0.063</b><br>(0.009)    | <b>-0.582</b><br>(0.024)   | <b>-0.069</b><br>(0.008)   |
| Panel C: (0.90, 0.90, 0.90) |                             |                            |                             |                              |                             |                             |                             |                            |                            |                             |                            |                            |
| Intercept                   | -124.840<br>(157.253)       | 161.731<br>(257.786)       | -72.439<br>(59.375)         | <b>-476.669</b><br>(76.346)  | <b>-202.209</b><br>(47.414) | <b>-82.551</b><br>(32.034)  | <b>-209.689</b><br>(24.663) | <b>491.864</b><br>(43.603) | 34.514<br>(22.358)         | <b>-148.541</b><br>(22.284) | <b>575.923</b><br>(59.029) | <b>-52.309</b><br>(19.866) |
| Wind Speed                  | <b>-1.137</b><br>(0.260)    | -0.369<br>(0.415)          | <b>-0.366</b><br>(0.151)    | <b>0.744</b><br>(0.121)      | <b>3.277</b><br>(0.074)     | <b>0.562</b><br>(0.053)     | <b>-0.441</b><br>(0.039)    | <b>0.704</b><br>(0.066)    | <b>-0.221</b><br>(0.035)   | -0.055<br>(0.034)           | <b>0.492</b><br>(0.091)    | <b>-0.278</b><br>(0.031)   |
| Temperature                 | <b>-0.725</b><br>(0.137)    | 0.178<br>(0.235)           | <b>-0.191</b><br>(0.053)    | <b>-1.844</b><br>(0.069)     | -0.009<br>(0.041)           | <b>-1.244</b><br>(0.029)    | <b>-0.261</b><br>(0.022)    | <b>0.494</b><br>(0.039)    | <b>-0.795</b><br>(0.019)   | <b>-0.149</b><br>(0.019)    | <b>-0.176</b><br>(0.049)   | <b>-0.428</b><br>(0.017)   |
| Pressure                    | 0.170<br>(0.152)            | -0.055<br>(0.247)          | 0.084<br>(0.057)            | <b>0.528</b><br>(0.074)      | <b>0.228</b><br>(0.046)     | <b>0.139</b><br>(0.031)     | <b>0.237</b><br>(0.024)     | <b>-0.423</b><br>(0.042)   | 0.009<br>(0.021)           | <b>0.165</b><br>(0.021)     | <b>-0.470</b><br>(0.057)   | <b>0.081</b><br>(0.019)    |
| Humidity                    | <b>-0.153</b><br>(0.064)    | <b>-0.551</b><br>(0.107)   | -0.016<br>(0.025)           | <b>-0.083</b><br>(0.033)     | <b>-0.212</b><br>(0.019)    | <b>-0.217</b><br>(0.013)    | <b>-0.056</b><br>(0.010)    | <b>-0.365</b><br>(0.018)   | <b>-0.105</b><br>(0.009)   | <b>-0.063</b><br>(0.009)    | <b>-0.582</b><br>(0.024)   | <b>-0.069</b><br>(0.008)   |



**Table 11** Estimated state-dependent correlation matrices for different quantile levels. Point estimates are displayed in boldface when significant at the standard 5% level

| States                      | $k = 1$                 |                          |                 | $k = 2$                  |                          |                 | $k = 3$                 |                          |                 | $k = 4$                 |                          |                 |
|-----------------------------|-------------------------|--------------------------|-----------------|--------------------------|--------------------------|-----------------|-------------------------|--------------------------|-----------------|-------------------------|--------------------------|-----------------|
|                             | PM <sub>25</sub>        | O <sub>3</sub>           | NO <sub>2</sub> | PM <sub>25</sub>         | O <sub>3</sub>           | NO <sub>2</sub> | PM <sub>25</sub>        | O <sub>3</sub>           | NO <sub>2</sub> | PM <sub>25</sub>        | O <sub>3</sub>           | NO <sub>2</sub> |
| Panel A: (0.50, 0.50, 0.50) |                         |                          |                 |                          |                          |                 |                         |                          |                 |                         |                          |                 |
| PM <sub>25</sub>            | 1                       |                          |                 | 1                        |                          |                 | 1                       |                          |                 |                         |                          |                 |
| O <sub>3</sub>              | 0.098<br>(0.062)        | 1                        |                 | -0.226<br>(0.154)        | 1                        |                 | 0.031<br>(0.057)        | 1                        |                 |                         |                          |                 |
| NO <sub>2</sub>             | <b>0.224</b><br>(0.063) | <b>-0.246</b><br>(0.061) | 1               | <b>0.512</b><br>(0.116)  | <b>-0.344</b><br>(0.137) | 1               | <b>0.391</b><br>(0.047) | -0.020<br>(0.056)        | 1               |                         |                          |                 |
| Panel B: (0.75, 0.75, 0.75) |                         |                          |                 |                          |                          |                 |                         |                          |                 |                         |                          |                 |
| PM <sub>25</sub>            | 1                       |                          |                 | 1                        |                          |                 | 1                       |                          |                 |                         |                          |                 |
| O <sub>3</sub>              | <b>0.206</b><br>(0.077) | 1                        |                 | -0.024<br>(0.095)        | 1                        |                 | <b>0.170</b><br>(0.056) | 1                        |                 |                         |                          |                 |
| NO <sub>2</sub>             | <b>0.199</b><br>(0.084) | <b>-0.422</b><br>(0.082) | 1               | <b>0.586</b><br>(0.061)  | -0.142<br>(0.097)        | 1               | <b>0.385</b><br>(0.050) | -0.021<br>(0.061)        | 1               |                         |                          |                 |
| Panel C: (0.90, 0.90, 0.90) |                         |                          |                 |                          |                          |                 |                         |                          |                 |                         |                          |                 |
| PM <sub>25</sub>            | 1                       |                          |                 | 1                        |                          |                 | 1                       |                          |                 | 1                       |                          |                 |
| O <sub>3</sub>              | 0.151<br>(0.193)        | 1                        |                 | <b>-0.361</b><br>(0.121) | 1                        |                 | 0.133<br>(0.076)        | 1                        |                 | -0.018<br>(0.094)       | 1                        |                 |
| NO <sub>2</sub>             | <b>0.385</b><br>(0.182) | <b>-0.591</b><br>(0.203) | 1               | <b>0.202</b><br>(0.077)  | <b>-0.749</b><br>(0.190) | 1               | 0.078<br>(0.077)        | <b>-0.230</b><br>(0.095) | 1               | <b>0.186</b><br>(0.089) | <b>-0.395</b><br>(0.114) | 1               |

PM<sub>25</sub>, O<sub>3</sub> and NO<sub>2</sub> in the wintertime is substantially higher than that in spring and summer.

### 6 Conclusion

The study of pollution exposure is at the heart of policy attention for health and economics welfare analysis. Motivated the necessity to develop sound policies for controlling air contaminants emissions, this paper extends the joint quantile regression of Petrella and Raponi (2019) by introducing a hidden semi-Markov quantile regression for the analysis of multiple pollutants time series. The proposed model allows to capture quantile-specific effects across the entire distribution of several outcomes in one step and infer cluster-specific covariate effects at various quantile levels of interest. In order to avoid making biased inference related to incorrect distributional assumptions about the SDs, we adopt the approach of Guédon (2003) where the latent sojourn densities are approximated by using nonparametric discrete distributions and estimated directly from the data. Using simulation exercises, the proposed approach reveals promising results in reducing state misclassification rates with respect to HMMs and identifying the correct number of hidden states. In the empirical application, we employ our methodology to jointly model daily PM<sub>25</sub>, O<sub>3</sub> and NO<sub>2</sub> concentrations recorded in Rieti (Italy) as a function of wind speed, temperature, pressure

and humidity. We find that seasonal changes in air pollutant concentrations are greatly affected by meteorological conditions, whose effects are generally amplified at the top end of the distribution of the responses. Moreover, the latent regimes capture seasonal variations in air pollution particles that characterize low and hazardous contamination levels.

This work could be extended further in the following directions. In particular, the method proposed could be positively applied to financial time series modeling. Indeed, financial returns often exhibit empirical characteristics, such as skewness, leptokurtosis, heteroscedasticity and clustering behavior over time, which are heavily influenced by hidden variables (e.g., the state of the market) during tranquil and crisis periods (Maruotti et al. 2021). In this context, time dependence can be further included in the regression model, allowing the quantile to vary over time according to an autoregressive process (Engle and Manganelli 2004).

The approach introduced might also be extended to other settings and data structures. Firstly, our QHSMM can be generalized to multivariate longitudinal data, extending the recent quantile mixed HMM in Merlo et al. (2022) by using more flexible sojourn distributions and, possibly, allowing the semi-Markov chain parameters to depend on observable covariates. Secondly, another extension would deal with the inclusion of spatial heterogeneity into the modeling approach to account for spatial-temporal variations of air pollutants from different monitoring sites. Lastly, while in the

application we focused on daily concentrations of three air pollutants, in high-dimensional settings with a larger number of response variables and/or number of hidden states, the introduced QHSMM can be easily over-parameterized. This often occurs because of the large number of unique parameters in the covariance matrices to be estimated, implying a loss in terms of interpretability as well as numerically ill-conditioned estimators. In these cases, following Maruotti et al. (2017), we may consider a class of parsimonious HSMMs by imposing a factor decomposition on the state-specific covariance matrices. Not only would this modeling strategy provide information about the dependence between pollutants, but also a clear interpretation of the latent association structure among them.

### Supplementary Materials

The Supplementary Materials include additional simulations that are used to support the results in the manuscript when the number of analyzed response variables  $p$  increases.

**Supplementary Information** The online version contains supplementary material available at <https://doi.org/10.1007/s11222-022-10130-1>.

**Acknowledgements** We would like to warmly thank the Associate Editor and two anonymous reviewers for their thoughtful comments and efforts towards improving our manuscript. This work was partially supported by the Finance Market Fund, Norway, project number 309218; “Statistical modelling and inference for (high-dimensional) financial data”.

### Appendix

**Proof of Proposition 1** Firstly, we note that to ensure identifiability of the MAL density in (5) it suffices to apply Proposition 1 of Petrella and Raponi (2019). Secondly, for a general HSMM, identifiability has been proven up to label switching (Leroux 1992). Thus, to ensure identifiability, all one needs to prove is the identifiability of the marginal mixtures (Dannemann et al. 2014), which in our case are represented by the finite mixtures of MAL distributions. Based on the work of Holzmann et al. (2006), Browne and McNicholas (2015) prove identifiability of finite mixtures of multivariate generalized hyperbolic distributions. Since the MAL in (5) is a limiting case of the multivariate generalized hyperbolic distribution (see (3) and (4) in Browne and McNicholas (2015) with  $\lambda = 1$ ,  $\psi = 2$  and  $\chi \rightarrow 0$ ), model identifiability follows by applying Corollary 2 of Browne and McNicholas (2015).

**Proof of Proposition 2** The E-step of the EM algorithm considers the conditional expectation of the complete log-

likelihood function in (8) given the observed data and the current parameter estimates  $\hat{\Phi}_\tau^{(r-1)}$ . At first, we recall that under the constraints imposed on  $\tilde{\xi}$  and  $\Lambda$ , the representation in (7) implies that:

$$\mathbf{Y} \mid \tilde{l} = \tilde{c} \sim \mathcal{N}_p(\boldsymbol{\mu} + \mathbf{D}\tilde{\xi}\tilde{c}, \tilde{c}\mathbf{D}\boldsymbol{\Sigma}\mathbf{D}), \quad \tilde{C} \sim \text{Exp}(1). \quad (35)$$

This means that the joint density function of  $\mathbf{Y}$  and  $\tilde{C}$  is:

$$f_{\mathbf{Y}, \tilde{C}}(\mathbf{y}, \tilde{c}) = \frac{\exp\left\{(\mathbf{y} - \boldsymbol{\mu})'\mathbf{D}^{-1}\boldsymbol{\Sigma}^{-1}\tilde{\xi}\right\}}{(2\pi)^{p/2} |\mathbf{D}\boldsymbol{\Sigma}\mathbf{D}|^{1/2}} \left(\tilde{c}^{-p/2} \exp\left\{-\frac{1}{2}\frac{\tilde{m}}{\tilde{c}} - \frac{1}{2}\tilde{c}(\tilde{d} + 2)\right\}\right). \quad (36)$$

By substituting (36) in (8) and taking the conditional expectation of the logarithm of (8), we obtain the expected complete log-likelihood function in (13).

To compute the conditional expectation of  $\tilde{c}_{tk}$  and  $\tilde{z}_{tk}$  in (13),  $\tilde{C}$  is treated as an additional latent variable. Using the joint distribution of  $\mathbf{Y}$  and  $\tilde{C}$  derived in (36) and the MAL density of  $\mathbf{Y}$  given in (5), we have that:

$$f_{\tilde{C}}(\tilde{C} \mid \mathbf{Y} = \mathbf{y}) = \frac{f_{\tilde{C}, \mathbf{Y}}(\tilde{c}, \mathbf{y})}{f_{\mathbf{Y}}(\mathbf{y})} = \frac{\tilde{c}^{-p/2} \left(\frac{2+\tilde{d}}{\tilde{m}}\right)^{v/2} \exp\left\{-\frac{\tilde{m}}{2\tilde{c}} - \frac{\tilde{c}(2+\tilde{d})}{2}\right\}}{2K_\nu\left(\sqrt{(2+\tilde{d})\tilde{m}}\right)}, \quad (37)$$

which corresponds to a Generalized Inverse Gaussian (GIG) distribution with parameters  $\nu, 2 + \tilde{d}, \tilde{m}_i$ , i.e.<sup>1</sup>

$$f_{\tilde{C}}(\tilde{C} \mid \mathbf{Y} = \mathbf{y}) \sim \text{GIG}\left(\nu, \tilde{d} + 2, \tilde{m}\right). \quad (38)$$

Then, it follows that

$$\mathbb{E}[\tilde{C} \mid \cdot] = \left(\frac{\hat{m}}{2 + \hat{d}}\right)^{\frac{1}{2}} \frac{K_{\nu+1}\left(\sqrt{(2 + \hat{d})\hat{m}}\right)}{K_\nu\left(\sqrt{(2 + \hat{d})\hat{m}}\right)} \quad (39)$$

and

$$\mathbb{E}[\tilde{C}^{-1} \mid \cdot] = \left(\frac{2 + \hat{d}}{\hat{m}}\right)^{\frac{1}{2}} \frac{K_{\nu+1}\left(\sqrt{(2 + \hat{d})\hat{m}}\right)}{K_\nu\left(\sqrt{(2 + \hat{d})\hat{m}}\right)} - \frac{2\nu}{\hat{m}}. \quad (40)$$

<sup>1</sup> The pdf of a GIG( $p, a, b$ ) distribution is defined as  $f_{GIG}(x; p, a, b) = \frac{(\frac{a}{b})^{p/2}}{2K_p(\sqrt{ab})} x^{p-1} e^{-\frac{1}{2}(ax+bx^{-1})}$ , with  $a > 0, b > 0$  and  $p \in \mathcal{R}$ .

Denoting the two conditional expectations in (39) and (40) by  $\hat{c}$  and  $\hat{z}$  respectively, concludes the proof.

**Proof of Proposition 3** Imposing the first order conditions on (13) with respect to each component of the set  $\Phi_\tau$ , gives the update estimates in (16), (17), (18) and (19). However, there is not closed formula solution to update the elements of the scale matrix  $D_j$ ; hence, the M-step update requires using numerical optimization techniques to maximize (13). A considerable disadvantage of this procedure is the necessary high computational effort which could be very time-consuming. For this reason, we utilize a simpler estimator for the scale parameters  $\delta_{jk}, k = 1, \dots, p$ , which follows directly from the fact that all marginals of the MAL distribution are univariate AL distributions (see Yu and Zhang 2005):

$$\hat{\delta}_{jk} = \frac{1}{T} \sum_{t=1}^T \hat{\gamma}_{tj} \rho_\tau(Y_t^{(k)} - \hat{\mu}_{tjk}), \tag{41}$$

where  $\hat{\mu}_{tjk}$  is the  $k$ -th element of the vector  $\hat{\mu}_{tj}$ .

## References

Adam, T., Langrock, R., Weiß, C.H.: Penalized estimation of flexible hidden Markov models for time series of counts. *Metron* **77**(2), 87–104 (2019)

Akaike, H.: Information theory and an extension of the maximum likelihood principle, in ‘Selected papers of Hirotugu Akaike’, Springer, pp. 199–213, (1998)

Barbu, V. S. and Limnios, N.: *Semi-Markov chains and hidden semi-Markov models toward applications: their use in reliability and DNA analysis*, Vol. 191, Springer Science & Business Media (2009)

Bartolucci, F., Farcomeni, A. and Pennoni, F.: *Latent Markov models for longitudinal data*, CRC Press, (2012)

Bassani, C., Vichi, F., Esposito, G., Montagnoli, M., Giusto, M., Ianniello, A.: Nitrogen dioxide reductions from satellite and surface observations during COVID-19 mitigation in Rome (Italy). *Environmental Science and Pollution Research* **28**(18), 22981–23004 (2021)

Bernardi, M., Gayraud, G., Petrella, L., et al.: Bayesian tail risk interdependence using quantile regression. *Bayesian Analysis* **10**(3), 553–603 (2015)

Biernacki, C., Celeux, G., Govaert, G.: Assessing a mixture model for clustering with the integrated completed likelihood. *IEEE Transactions on Pattern Analysis and Machine Intelligence* **22**(7), 719–725 (2000)

Browne, R.P., McNicholas, P.D.: A mixture of generalized hyperbolic distributions. *Canadian Journal of Statistics* **43**(2), 176–198 (2015)

Bulla, J., Bulla, I.: Stylized facts of financial time series and hidden semi-Markov models. *Computational Statistics & Data Analysis* **51**(4), 2192–2209 (2006)

Bulla, J., Bulla, I., Nenadić, O.: hsmm-an R package for analyzing hidden semi-Markov models. *Computational Statistics & Data Analysis* **54**(3), 611–619 (2010)

Cappé, O., Moulines, E. and Rydén, T. (2006), *Inference in hidden Markov models*, Springer Science & Business Media

Charlier, I., Paindaveine, D., Saracco, J.: Multiple-output quantile regression through optimal quantization. *Scandinavian Journal of Statistics* **47**(1), 250–278 (2020)

Chavas, J.-P.: On multivariate quantile regression analysis. *Statistical Methods & Applications* **27**(3), 365–384 (2018)

Dannemann, J., Holzmann, H., Leister, A.: Semiparametric hidden Markov models: identifiability and estimation. *Wiley Interdisciplinary Reviews: Computational Statistics* **6**(6), 418–425 (2014)

Dempster, A. P., Laird, N.M. and Rubin, D.B.: ‘Maximum likelihood from incomplete data via the EM algorithm’, *Journal of the Royal Statistical Society Series B (Methodological)* pp. 1–38 (1977)

El Ghouch, A., Genton, M.G.: Local polynomial quantile regression with parametric features. *J. Am. Stat. Assoc.* **104**(488), 1416–1429 (2009)

Engle, R.F., Manganelli, S.: CAViaR: conditional autoregressive value at risk by regression quantiles. *J. Bus. Econ. Stat.* **22**(4), 367–381 (2004)

Ephraim, Y., Merhav, N.: Hidden Markov processes. *IEEE Trans. Inf. Theor.* **48**(6), 1518–1569 (2002)

Farcomeni, A.: Quantile regression for longitudinal data based on latent Markov subject-specific parameters. *Stat. Comput.* **22**(1), 141–152 (2012)

Geraci, M.: Modelling and estimation of nonlinear quantile regression with clustered data. *Comput. Stat. Data Anal.* **136**, 30–46 (2019)

Guédon, Y.: Estimating hidden semi-Markov chains from discrete sequences. *J. Comput. Graph. Stat.* **12**(3), 604–639 (2003)

Hamilton, J. D.: ‘A new approach to the economic analysis of nonstationary time series and the business cycle’, *Econometrica: Journal of the Econometric Society* pp. 357–384, (1989)

Holzmann, H., Munk, A., Gneiting, T.: Identifiability of finite mixtures of elliptical distributions. *Scandinav. J. Stat.* **33**(4), 753–763 (2006)

Hubert, L., Arabie, P.: Comparing partitions. *J. Classif.* **2**(1), 193–218 (1985)

Koenker, R.: *Quantile regression*. Cambridge University Press (2005)

Koenker, R. and Bassett, G.: ‘Regression Quantiles’, *Econometrica: Journal of the Econometric Society* **46**(1), 33–50, (1978)

Koenker, R., Chernozhukov, V., He, X. and Peng, L.: *Handbook of quantile regression*, CRC press, (2017)

Kong, L. and Mizera, I.: ‘Quantile tomography: using quantiles with multivariate data’, *Statistica Sinica* pp. 1589–1610, (2012)

Kotz, S., Kozubowski, T. and Podgorski, K.: *The Laplace distribution and generalizations: a revisit with applications to communications, economics, engineering, and finance*, Springer Science & Business Media, (2012)

Langrock, R., Kneib, T., Sohn, A., DeRuiter, S.L.: Nonparametric inference in hidden Markov models using P-splines. *Biometrics* **71**(2), 520–528 (2015)

Langrock, R., Zucchini, W.: Hidden Markov models with arbitrary state dwell-time distributions. *Comput. Stat. Data Anal.* **55**(1), 715–724 (2011)

Leroux, B.G.: Maximum-likelihood estimation for hidden Markov models. *Stochastic Process. Appl.* **40**(1), 127–143 (1992)

Levinson, S.E., Rabiner, L.R., Sondhi, M.M.: An introduction to the application of the theory of probabilistic functions of a Markov process to automatic speech recognition. *Bell Syst. Tech. J.* **62**(4), 1035–1074 (1983)

Luo, Y., Lian, H., Tian, M.: Bayesian quantile regression for longitudinal data models. *J. Stat. Comput. Simul.* **82**(11), 1635–1649 (2012)

MacDonald, I. L. and Zucchini, W.: *Hidden Markov and other models for discrete-valued time series*, Vol. 110, CRC Press (1997)

Marino, M.F., Tzavidis, N., Alfò, M.: Mixed hidden Markov quantile regression models for longitudinal data with possibly incomplete sequences. *Stat. Methods Med. Res.* **27**(7), 2231–2246 (2018)

Maruotti, A.: Mixed hidden Markov models for longitudinal data: an overview. *Int. Stat. Rev.* **79**(3), 427–454 (2011)

- Maruotti, A., Bulla, J., Lagona, F., Picone, M. and Martella, F.: 'Dynamic mixtures of factor analyzers to characterize multivariate air pollutant exposures', *The Annals of Applied Statistics* pp. 1617–1648, (2017)
- Maruotti, A., Petrella, L., Sposito, L.: Hidden semi-Markov-switching quantile regression for time series. *Compu. Stati. Data Anal.* **159**, 107208 (2021)
- Maruotti, A., Punzo, A.: Initialization of hidden Markov and semi-Markov models: a critical evaluation of several strategies. *Int. Stat. Rev.* **89**(3), 447–480 (2021)
- Maruotti, A., Punzo, A., Bagnato, L.: Hidden Markov and semi-Markov models with multivariate leptokurtic-normal components for robust modeling of daily returns series. *J. Financ. Econom.* **171**(1), 91–117 (2019)
- Merlo, L., Petrella, L., Raponi, V.: Forecasting VaR and ES using a joint quantile regression and its implications in portfolio allocation. *J. Bank. Finan.* **133**, 106248 (2021)
- Merlo, L., Petrella, L., Salvati, N., Tzavidis, N.: Marginal M-quantile regression for multivariate dependent data. *Comput. Stat. Data Anal.* **173**, 107500 (2022)
- Merlo, L., Petrella, L., Tzavidis, N.: 'Quantile mixed hidden Markov models for multivariate longitudinal data: an application to children's Strengths and Difficulties Questionnaire scores', *Journal of the Royal Statistical Society. Ser. C Appl. Stat.* **71**(2), 417–448 (2022)
- O'Connell, J., Højsgaard, S., et al.: Hidden semi Markov models for multiple observation sequences: the mhsmm package for R. *J. Stat. Softw.* **39**(4), 1–22 (2011)
- Petrella, L., Raponi, V.: Joint estimation of conditional quantiles in multivariate linear regression models with an application to financial distress. *J. Multivar. Anal.* **173**, 70–84 (2019)
- Pohle, J., Adam, T., Beumer, L.T.: Flexible estimation of the state dwell-time distribution in hidden semi-Markov models. *Comput. Stat. Data Anal.* **172**, 107479 (2022)
- Pohle, J., Langrock, R., van Beest, F.M., Schmidt, N.M.: Selecting the number of states in hidden Markov models: pragmatic solutions illustrated using animal movement. *J. Agric Biol. Environ. Stat.* **22**(3), 270–293 (2017)
- Putaud, J.-P., Pozzoli, L., Pisoni, E., Martins Dos Santos, S., Lagler, F., Lanzani, G., Dal Santo, U., Colette, A.: Impacts of the COVID-19 lockdown on air pollution at regional and urban background sites in northern Italy. *Atmosp. Chem. Phys.* **21**(10), 7597–7609 (2021)
- Sansom, J. and Thomson, P. (2001), 'Fitting hidden semi-Markov models to breakpoint rainfall data', *Journal of Applied Probability* **38**(A), 142–157
- Schwarz, G., et al.: Estimating the dimension of a model. *Anna. Stat.* **6**(2), 461–464 (1978)
- Serfling, R.: Quantile functions for multivariate analysis: approaches and applications. *Statist. Neerlandica* **56**(2), 214–232 (2002)
- Stolfi, P., Bernardi, M. and Petrella, L.: 'The sparse method of simulated quantiles: An application to portfolio optimization', *Statistica Neerlandica* (2018),
- Visser, I., Raijmakers, M.E., Molenaar, P.C.: Confidence intervals for hidden Markov model parameters. *Br. J. Math. Statist. Psychol.* **53**(2), 317–327 (2000)
- Ye, W., Zhu, Y., Wu, Y. and Miao, B.: 'Markov regime-switching quantile regression models and financial contagion detection', *Insurance: Mathematics and Economics* **67**, 21–26, (2016)
- Yu, K., Moyeed, R.A.: Bayesian quantile regression. *Stat. Probab. Lett.* **54**(4), 437–447 (2001)
- Yu, K., Zhang, J.: A three-parameter asymmetric Laplace distribution and its extension. *Commun. Statist. Theory Methods* **34**(9–10), 1867–1879 (2005)
- Yu, S.-Z.: Hidden Semi-Markov models: theory, algorithms and applications. Morgan Kaufmann (2015)
- Zucchini, W., MacDonald, I. L. and Langrock, R.: *Hidden Markov models for time series: an introduction using R*, Chapman and Hall CRC (2016)

**Publisher's Note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Springer Nature or its licensor holds exclusive rights to this article under a publishing agreement with the author(s) or other rightsholder(s); author self-archiving of the accepted manuscript version of this article is solely governed by the terms of such publishing agreement and applicable law.