

Bayesian hidden Markov modelling using circular-linear general projected normal distribution

Original

Bayesian hidden Markov modelling using circular-linear general projected normal distribution / Mastrantonio, Gianluca; Maruotti, Antonello; Jona Lasinio, Giovanna. - In: ENVIRONMETRICS. - ISSN 1180-4009. - 26:2(2015), pp. 145-158. [10.1002/env.2326]

Availability:

This version is available at: 11583/2664908 since: 2020-01-30T11:33:49Z

Publisher:

John Wiley and Sons

Published

DOI:10.1002/env.2326

Terms of use:

This article is made available under terms and conditions as specified in the corresponding bibliographic description in the repository

Publisher copyright

(Article begins on next page)

Bayesian Hidden Markov Modelling Using Circular-Linear General Projected Normal Distribution

Gianluca Mastrantonio *

Department of Economics, University of Romatre

and

Antonello Maruotti †

Southampton Statistical Science Research Institute, University of Southampton

and

Giovanna Jona Lasinio ‡

Department of Statistical Science, Sapienza University of Rome

Abstract

We introduce a multivariate hidden Markov model to jointly cluster time-series observations with different support, i.e. circular and linear. Relying on the general projected normal distribution, our approach allows for bimodal and/or skewed cluster-specific distributions for the circular variable. Furthermore, we relax the independence assumption between the circular and linear components observed at the same time. Such an assumption is generally used to alleviate the computational burden involved in the parameter estimation step, but it is hard to justify in empirical applications. We carry out a simulation study using different data-generation schemes to investigate model behavior, focusing on well recovering the hidden structure. Finally, the model is used to fit a real data example on a bivariate time series of wind speed and direction.

1 INTRODUCTION

Hidden Markov models (HMMs) have become more frequently used to provide a natural and flexible framework for univariate and multivariate time-dependent data (e.g. time-series, longitudinal data). They are a class of mixture models in which the data-generation distribution depends on the state of an underlying and unobserved Markov process. Hidden Markov modelling has been used as a statistical tool for density estimation (Langrock *et al.*, 2014; Dannemann, 2012), supervised and unsupervised classification (Lagona and Picone, 2012; Alfò and Maruotti, 2010; Frühwirth Schnatter, 2006) and a wide range of empirical problems in environmetrics (Martinez-Zarzoso and Maruotti, 2013; Langrock *et al.*, 2012a), medicine (Langrock *et al.*, 2013; Lagona *et al.*, 2014), education (Bartolucci *et al.*, 2011). For a comprehensive introduction to fundamental theory of HMMs encountered in practice, see the review papers of Bartolucci *et al.* (2014), Maruotti (2011) and monographs by Bartolucci *et al.* (2012), Zucchini and MacDonald (2009) and Cappé *et al.* (2005).

The literature on multivariate hidden Markov modelling is dominated by Gaussian HMMs (Spezia, 2010; Bartolucci and Farcomeni, 2010; Geweke and Amisano, 2011). Modelling multivariate time series with non-normal components of mixed-type is challenging. The joint distribution of multivariate (mixed-type) data is usually specified as a mixture having products of univariate distributions as components (see e.g. Lagona *et al.*, 2011; Lagona and Picone, 2011; Zhang *et al.*, 2010). Bartolucci and Farcomeni (2009) is a notable exception. In other words, random variables are assumed conditionally independent given the latent structure. Although conditional independence facilitates parameters estimation, it is a too restrictive assumption in many empirical applications and may not properly accommodate for the complex shape of multivariate distributions (Baudry *et al.*, 2010). Moreover,

*gianluca.mastrantonio@uniroma3.it Corresponding author

†a.maruotti@soton.ac.uk

‡giovanna.jonalasinio@uniroma1.it

an unnecessary number of latent states is often needed to obtain reasonable fit, at the price of an increased computational burden and difficulties in interpreting results, as shown in the simulation study section.

In this paper, we propose a bivariate distribution for circular-linear time-series in a HMM framework. We accommodate for nonstandard features of data including correlation in time and across variables, mixed supports (circular and linear) of the data, the special nature of circular measurements, and the occurrence of missing values. We relax the conditional independence assumption between circular and linear variables by taking a fully parametric approach.

This is not the first attempt to jointly modelling circular and linear variables in a HMM framework. Bulla *et al.* (2012) introduced a latent-class approach to the analysis of multivariate mixed-type data by assuming that circular and linear variables are conditionally independent given the states visited by a latent Markov chain while Kato *et al.* (2008) propose a hyper-cylindrical distribution. The latter is problematic (in the HMM setting in particular), because little is known about efficient estimation procedures and identifiability issues under hyper-cylindrical parametric models. In addition, mixtures of hyper-toroidal densities would group data according to clusters of difficult interpretation, without necessarily improving the fit of the model.

We introduce a flexible structure, relying on the general projected normal distribution (Wang and Gelfand, 2014), to model circular measurements, and extending Bulla *et al.* (2012) to a more general setting, allowing for (conditional) correlation between circular and linear variables. We treat the circular response as projection onto the unit circle of a bivariate variable and define the joint circular-linear distribution through the specification of a multivariate model in a multivariate linear setting, extending Wang *et al.* (2014) to a clustering framework.

The resulting hidden Markov model parameters are estimated in a Bayesian framework. We provide details on how to fit the model by using MCMC methods, and we point out possible drawbacks in the implementation of the algorithm. Advantages of the Bayesian approach, with respect to the EM algorithm, include a convenient framework to simultaneously account for several data features, adjust for identifiability issues and produce natural measures of uncertainty for model parameters. For a general discussion see e.g. (Rydén and Titterton, 1998; Rydén, 2008; Yildirim *et al.*, 2014).

We illustrate the proposal by a large-scale simulation study in order to investigate the empirical behaviour of the proposed approach with respect to several factors, such as the number of observed times, the association structure between the circular and linear variables and the fuzziness of the classification. We evaluate model performance in recovering the true model structure, we compare several models on the basis of their ability to accurately estimate the vectors of state-dependent parameters and hidden parameters. Finally, we test the proposal by analysing time series of semi-hourly wind directions and speeds, recorded in the period 12/12/2009, 12/1/2010 by the buoy of Ancona, located in the Adriatic Sea at about 30 km from the coast.

The rest of the paper is organized as follows. In Section 2, we briefly review relevant aspects necessary for the introduction of our approach and outline some results about the projected normal distribution. Section 3 discusses the specification of the circular-linear general projected normal hidden Markov model and provides Bayesian inference. Computational details and parameters estimation are discussed as well. Section 4 presents a large-scale simulation study. In Section 5, the application of the proposed methodology is illustrated through a real-world data set. Some concluding remarks are given in Section 6.

2 PRELIMINARIES

Circular data are a particular class of directional data, specifically, they are directions in two dimensions. To analyze circular data is challenging because usual statistics, which have been developed for linear data (for example the mean and variance), will not be meaningful and will be misleading when applied to directional data without taking into account the particular definition of the domain. There are many ways to define distributions in a circular domain, see the book of Mardia and Jupp (1999) for a comprehensive overview. The one we used in this paper is to radially project onto the circle a probability distribution originally defined on the plane. Let $\mathbf{Z} = [Z_1, Z_2]'$ be a 2-dimensional random

vector such that $\Pr(\mathbf{Z} = \mathbf{0}) = 0$. Then, its radial projection $\mathbf{W} = \begin{bmatrix} W_1 \\ W_2 \end{bmatrix} = \frac{\mathbf{Z}}{\|\mathbf{Z}\|}$ is a random vector on the unit circle, which can be converted to a random angle X relative to some direction treated as 0 via the transformation $X = \arctan^* \frac{W_2}{W_1} = \arctan^* \frac{Z_2}{Z_1} \in [0, 2\pi)$, where the function \arctan^* is a quadrant specific inverse of the tangent function, sometimes called *atan2*, that takes into account the signs of W_1 and W_2 to identify the right quadrant of X ; for a formal definition see Jammalamadaka and SenGupta (2001), pag. 13. Note that $\mathbf{W} = \begin{bmatrix} \cos X \\ \sin X \end{bmatrix}$ and let $R = \|\mathbf{Z}\|$ the following relation

$$\text{holds: } \begin{bmatrix} Z_1 \\ Z_2 \end{bmatrix} = \begin{bmatrix} R \cos X \\ R \sin X \end{bmatrix} = R\mathbf{W}.$$

By assuming $\mathbf{Z} \sim N_2(\cdot|\boldsymbol{\mu}, \boldsymbol{\Sigma})$, with $\boldsymbol{\Sigma} = \begin{bmatrix} \sigma_1^2 & \sigma_1\sigma_2\rho \\ \sigma_1\sigma_2\rho & \sigma_2^2 \end{bmatrix}$ and $\boldsymbol{\mu} = \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}$, X is said to have a 2-dimensional projected normal distribution, denoted by $PN_2(\cdot|\boldsymbol{\mu}, \boldsymbol{\Sigma})$. Since the distribution of X does not change if we multiply \mathbf{Z} for a positive constant $c > 0$, for identifiability purposes, following Wang and Gelfand (2012), σ_2^2 is set to be 1. The projected normal distribution is specified as a four parameters distribution: $PN_2(\cdot|\mu_1, \mu_2, \sigma_1^2, \rho)$.

We provide some examples to illustrate the flexibility of the projected normal (PN) distribution. The PN density can be symmetric, asymmetric or possibly bimodal and apart from some special case, the interpretation of the parameters can be difficult. The number of modes and the shape depend on the value of the all parameters and different sets of parameters can give really similar shapes. As a general comment we highlight that μ_1 and μ_2 are the means of the two Cartesian coordinates z_1 and z_2 and are respectively connected to the cosine and sine of the circular variable. By fixing $\sigma_1^2 = 1$ and $\rho = 0$, the resulting distribution is unimodal and symmetric and if $\mu_1 = \mu_2 = 0$ the distribution becomes a circular uniform. Departure from zero for the two means, in the case of identity covariance matrix, creates one mode in the trigonometric quadrant with the same sign of the means, e.g. if $\mu_1 > 0$ and $\mu_2 < 0$ then the mode is in the quadrant with positive cosine and negative sine; higher values of a mean attract the mode to its correspondent axis. By allowing the ρ parameter to vary, we obtain very flexible shapes. The resulting distribution shows asymmetry with more mass of probability near the axis with the highest μ . By increasing $|\rho|$, bimodality is detected, Figure 1. Moreover, for $\sigma_1^2 < 1$ the modes are closer to the sine axis; while for $\sigma_1^2 > 1$ the modes are closer to the cosine axis.

3 THE CIRCULAR-LINEAR GENERAL PROJECTED NORMAL HIDDEN MARKOV MODEL

3.1 The Model

Let $\mathcal{T} = \{0, 1, \dots, T-1, T\}$, in this paper we consider a bivariate time series $[\mathbf{x}, \mathbf{y}] = \{[x_t, y_t]; t \in \mathcal{T} \setminus \{0\}\}$ with circular, x_t , and linear, y_t , components. Our aim is to jointly classify $[x_t, y_t]$ in K classes, generally called regimes or states, with a HMM-based classifier.

Let $\pi_{k,h}$ indicates the probability to move from state k to state h and let ξ_{tk} be an indicator variable such that if we are in state k on time t it is 1, otherwise is 0. Then $f(\xi_{th} = 1 | \xi_{t-1k} = 1) = \pi_{k,h}$ and we set $f(\xi_{0k}) = \pi_k$. We indicate with

$$\boldsymbol{\pi} = \begin{bmatrix} \pi_{1,1} & \pi_{1,2} & \cdots & \pi_{1,K} \\ \pi_{2,1} & \pi_{2,2} & \cdots & \pi_{2,K} \\ \cdots & \cdots & \cdots & \cdots \\ \pi_{K,1} & \pi_{K,2} & \cdots & \pi_{K,K} \end{bmatrix}, \quad \sum_{h=1}^K \pi_{k,h} = 1, \quad k = 1, 2, \dots, K,$$

the transition matrix that governs the evolution of the Markov chain, $\boldsymbol{\pi}_0 = [\pi_1, \pi_2, \dots, \pi_K]'$ and $\boldsymbol{\xi} = [\boldsymbol{\xi}_0, \boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_T]'$ where $\boldsymbol{\xi}_t = [\xi_{t1}, \xi_{t2}, \dots, \xi_{tK}]'$.

Let $n_{k,h} = \sum_{t=1}^T \xi_{t-1k} \xi_{th}$ be the number of times we move from state k to state h , the joint density of the vector of states is $f(\boldsymbol{\xi} | \boldsymbol{\pi}, \boldsymbol{\pi}_0) = \prod_{k=1}^K \pi_k^{\xi_{0k}} \prod_{k=1}^K \prod_{h=1}^K \pi_{k,h}^{n_{k,h}}$. In the literature on HMM for circular-linear variables, see for example Bulla *et al.* (2012) and Holzmamann *et al.* (2006), it is generally assumed that conditioning to the latent vector $\boldsymbol{\xi}$, the pairs $[x_t, y_t]$ and $[x_g, y_g]$ are independent if $g \neq t$ and at the

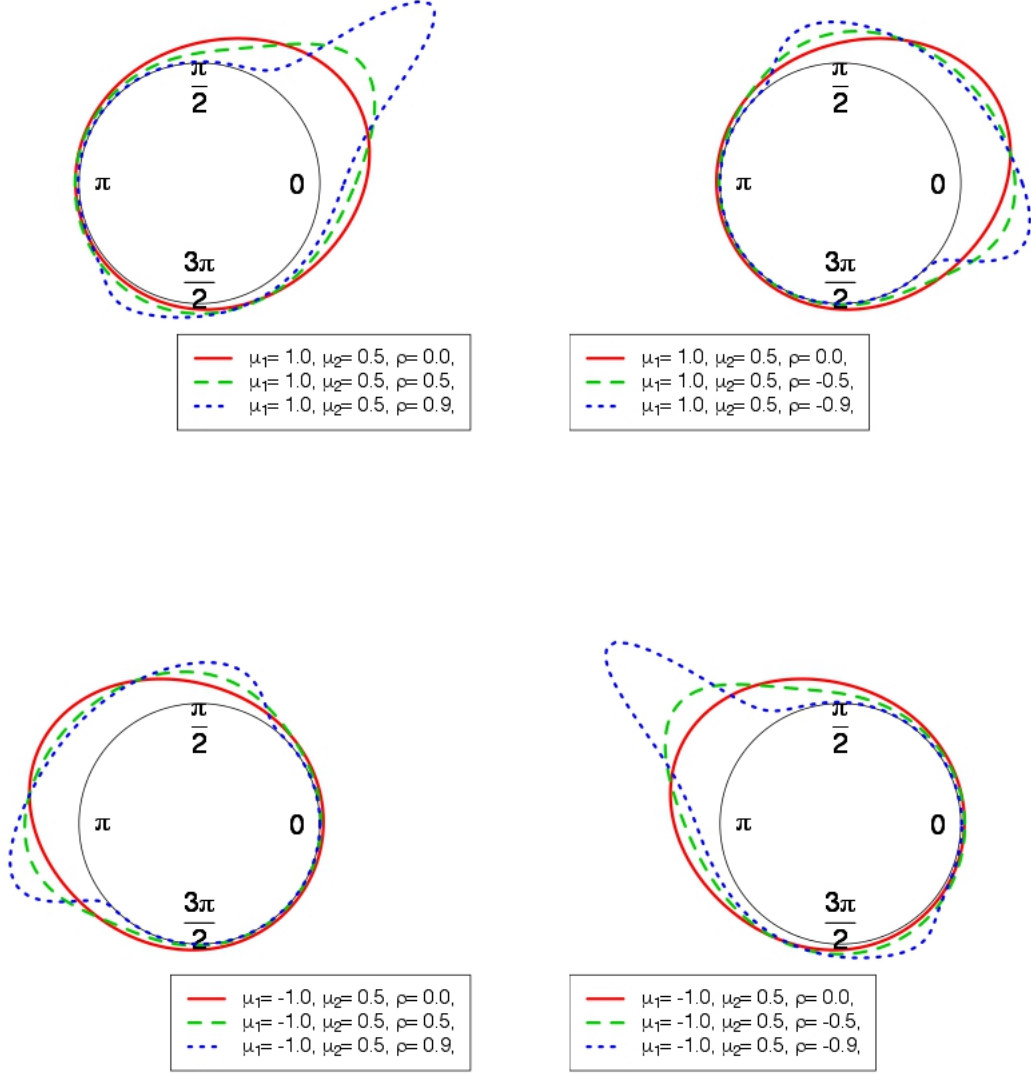


Figure 1: Shape of the projected normal distribution for $\sigma^2 = 1$ and different values of ρ , μ_1 and μ_2

same time $x_t \perp y_t$. As a result, the conditional distribution of the observed process, given the latent process, takes the form of a product density, say $f(\mathbf{x}, \mathbf{y}|\boldsymbol{\xi}) = \prod_{k=1}^K \prod_{t=1}^T [f(x_t|\xi_{tk} = 1)f(y_t|\xi_{tk} = 1)]^{\xi_{tk}}$. We maintain the so-called conditional independence property: given the hidden state at time t , the distribution of the observation at this time is fully determined, but we relax the assumption on independence between the circular and linear variables observed at the same time. Thus, we get a multivariate conditional distribution $f(\mathbf{x}, \mathbf{y}|\boldsymbol{\xi}) = \prod_{k=1}^K \prod_{t=1}^T f(x_t, y_t|\xi_{tk} = 1)^{\xi_{tk}}$.

Let $\mathbf{Z}_t|\xi_{tk} = 1 \sim N(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$, with $\mathbf{Z}_t = \begin{bmatrix} Z_{t1} \\ Z_{t2} \end{bmatrix}$, $\boldsymbol{\mu}_k = \begin{bmatrix} \mu_{k1} \\ \mu_{k2} \end{bmatrix}$, $\boldsymbol{\Sigma}_k = \begin{bmatrix} \sigma_{k1}^2 & \sigma_{k1}\rho_k \\ \sigma_{k1}\rho_k & 1 \end{bmatrix}$ and let $R_t = \|\mathbf{Z}_t\|$. We define X_t as the radial projection of \mathbf{Z}_t : $X_t = \arctan^* \frac{Z_{t1}}{Z_{t2}}$ and then $X_t|\xi_{tk} = 1 \sim PN_2(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$. We can write easily the joint density of $[X_t, R_t]$ that is the density that arises by a variable transformation from the bivariate normal Z_t to its polar system representation. Let $\phi_h(\zeta|\mathbf{M}, \mathbf{V})$ be the probability density function of a h -variate normal distribution with mean \mathbf{M} and covariance matrix \mathbf{V} evaluated in ζ , then

$$f(x_t, r_t|\xi_{tk} = 1) = \phi_2(r_t \mathbf{w}_t | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) r_t \quad (1)$$

We built the (conditional) joint density $f(x_t, y_t | \xi_{tk} = 1) = f(y_t | x_t, \xi_{tk} = 1) f(x_t | \xi_{tk} = 1)$ as a marginalization over the latent variable R_t : $f(x_t, y_t | \xi_{tk} = 1) = \int_{r_t} f(y_t | x_t, r_t, \xi_{tk} = 1) f(x_t, r_t | \xi_{tk} = 1) dr_t$, where $f(x_t, r_t | \xi_{tk} = 1)$ is specified in equation (1) and $y_t | x_t, r_t, \xi_{tk} = 1$ is defined through a circular-linear regression. However, there is not an obvious and standard way to formalize the relations between circular and linear variables. Jammalamadaka and SenGupta (2001) propose a flexible approach using trigonometric polynomials while Mardia (1976) and Johnson and Wehrly (1978) proposed a regression where the covariates are the sine and cosine components of the circular variable. Here, following Wang *et al.* (2014), we specify the relation as $y_t = \gamma_{k0} + \gamma_{k1} r_t \cos x_t + \gamma_{k2} r_t \sin x_t + \epsilon_{tk}$, with $\epsilon_{tk} \sim N(0, \sigma_{ky}^2)$. Thus, $y_t | x_t, r_t, \xi_{tk} = 1$ is distributed as a normal variable with mean $\gamma_{k0} + \gamma_{k1} r_t \cos x_t + \gamma_{k2} r_t \sin x_t$ and variance σ_{ky}^2 . Note that the regression can be seen as a linear regression between y_t and the inline variables $r_t \cos x_t$ and $r_t \sin x_t$. This type of representation gives more flexibility to the circular linear regression than the ones proposed by Mardia (1976) and Johnson and Wehrly (1978). Notice that with the r_t variable, for a given value of the circular variable at different time point, say $x_t = x_{t'}, t \neq t'$, the relation between x_t and y_t and $x_{t'}$ and $y_{t'}$ can be different as it depends on the realization of the non observed variable r_t .

Then we have that

$$f(x_t, y_t | \xi_{tk} = 1) = \frac{\phi_1(y_t | \gamma_{k0}, \sigma_{ky}^2) \phi_2(\boldsymbol{\mu}_k | \mathbf{0}_2, \boldsymbol{\Sigma}_k) \left[m_{tk} \Phi\left(\frac{m_{tk}}{\sqrt{v_{tk}}}\right) + \phi_1(m_{tk} | 0, v_{tk}) \right]}{\phi_1(m_{tk} | 0, v_{tk})}, \quad (2)$$

where Φ is the cumulative density function of a standard normal distribution, $\mathbf{w} = \begin{bmatrix} \cos x_t \\ \sin x_t \end{bmatrix}$, $v_{tk} = \left[\frac{c_{tk}^2}{\sigma_{ky}^2} + \mathbf{w}'_t \boldsymbol{\Sigma}_k^{-1} \mathbf{w}_t \right]^{-1}$, $m_{tk} = v_{tk} \left[\frac{c_{tk}(y_t - \gamma_{k0})}{\sigma_{ky}^2} + \mathbf{w}'_t \boldsymbol{\Sigma}_k^{-1} \boldsymbol{\mu}_k \right]$ and $c_{tk} = \mathbf{w}'_t \begin{bmatrix} \gamma_{k1} \\ \gamma_{k2} \end{bmatrix}$. The Circular Linear General projected normal (CL-GPN) distribution with parameters $[\mu_{k1}, \mu_{k2}, \sigma_{k1}^2, \rho_k, \gamma_{k0}, \gamma_{k1}, \gamma_{k2}, \sigma_{ky}^2]'$ is thus defined in (2). In this setting the parameter γ_{k1} and γ_{k2} govern the dependency between the two variables (linear and circular), γ_{k1} is connected to the correlation between the linear variable and the cosine of the circular, γ_{k1} is connected to the correlation between the linear variable and the sine of the circular.

Wang *et al.* (2014) and Wang and Gelfand (2012) argue that working with the projected normal density is not easy and its form is practically intractable (to see the closed form of the PN density see Wang and Gelfand (2012)). Since the CL-GPN is based on the PN, is itself an intractable distribution and the implementation of the MCMC algorithm can be difficult. However the introduction of r_t is of practical use as it simplifies the implementation of the MCMC algorithm, see Section 3.3.

3.2 Posterior Inference

Let $\boldsymbol{\Psi}$ be the vector of all the parameters of the CL-GPN in all the K regimes, we have the following posterior distribution

$$f(\boldsymbol{\pi}, \boldsymbol{\xi}, \boldsymbol{\pi}_0, \boldsymbol{\Psi}, \mathbf{r} | \mathbf{x}, \mathbf{y}) = \frac{f(\mathbf{r}, \mathbf{x}, \mathbf{y} | \boldsymbol{\Psi}, \boldsymbol{\xi}) f(\boldsymbol{\xi}_{-0} | \boldsymbol{\xi}_0, \boldsymbol{\pi}) f(\boldsymbol{\pi}) f(\boldsymbol{\xi}_0 | \boldsymbol{\pi}_0) f(\boldsymbol{\pi}_0) f(\boldsymbol{\Psi})}{f(\mathbf{x}, \mathbf{y})}$$

where $\mathbf{r} = [r_1, \dots, r_t]'$ and $f(\mathbf{x}, \mathbf{y}, \mathbf{r} | \boldsymbol{\xi}) = \prod_{k=1}^K \prod_{t=1}^T f(x_t, r_t, y_t | \xi_{tk} = 1)^{\xi_{tk}}$. As prior distribution we assume: $\mu_{ki} \sim N(\cdot, \cdot)$, $\sigma_{k1}^2 \sim IG(\cdot, \cdot)$, $\rho_k \sim N(\cdot, \cdot) I(-1, 1)$, $\sigma_{ky}^2 \sim IG(\cdot, \cdot)$, $\gamma_{kj} \sim N(\cdot, \cdot)$ for $k = 1, \dots, K$, $i = 1, 2$, $j = 1, 2, 3$, where $IG(\cdot, \cdot)$ indicates the Inverse Gamma distribution, $\boldsymbol{\pi}_0 \sim Dir(\cdot)$ and $\boldsymbol{\pi}_{k,\cdot} \sim Dir(\cdot)$ where $Dir(\cdot)$ indicates the Dirichlet distribution and $\boldsymbol{\pi}_{k,\cdot}$ is the k^{th} row of $\boldsymbol{\pi}$: we assume $\boldsymbol{\pi}_k \perp \boldsymbol{\pi}_{k'}$ if $k \neq k'$. The prior specification allows us to marginalized over $\boldsymbol{\pi}$ and $\boldsymbol{\pi}_0$ reducing of K^2 the number of parameters to simulate and leads to a more efficient and stable algorithm (see Banerjee *et al.* (2004) and Section 3.3). Note that we can always sample from $f(\boldsymbol{\pi}_{k,\cdot} | \mathbf{x}, \mathbf{y}) = \sum_{\boldsymbol{\xi}} f(\boldsymbol{\pi}_{k,\cdot} | \boldsymbol{\xi}) f(\boldsymbol{\xi} | \mathbf{x}, \mathbf{y})$ and $f(\boldsymbol{\pi}_0 | \mathbf{x}, \mathbf{y}) = \sum_{\boldsymbol{\xi}} f(\boldsymbol{\pi}_0 | \boldsymbol{\xi}) f(\boldsymbol{\xi} | \mathbf{x}, \mathbf{y})$ given the set of B posterior samples $\{\boldsymbol{\xi}^b\}_{b=1}^B$ of $\boldsymbol{\xi}$, with an MCMC integration. For each sample $\boldsymbol{\xi}^b$ we draw a sample from $\boldsymbol{\pi}_{k,\cdot}^b | \boldsymbol{\xi}^b \sim Dir(\beta + \sum_{t=2}^T \xi_{t-1k}^b \xi_{t1}^b, \dots, \beta + \sum_{t=2}^T \xi_{t-1k}^b \xi_{tK}^b)$ and one from $\boldsymbol{\pi}_0^b | \boldsymbol{\xi}^b \sim Dir(\beta + \xi_{01}^b, \dots, \beta + \xi_{0K}^b)$. The sets $\{\boldsymbol{\pi}_{k,\cdot}^b\}_{b=1}^B$ and

$\{\pi_0^b\}_{b=1}^B$ are drawn from their respective marginal posterior distributions. The posterior distribution we will work with is then

$$f(\boldsymbol{\xi}, \boldsymbol{\Psi}, \mathbf{r}|\mathbf{x}, \mathbf{y}) = \frac{f(\mathbf{r}, \mathbf{x}, \mathbf{y}|\boldsymbol{\Psi}, \boldsymbol{\xi}) f(\boldsymbol{\Psi}) f(\boldsymbol{\xi}_{-0}|\boldsymbol{\xi}_0) f(\boldsymbol{\xi}_0)}{f(\mathbf{x}, \mathbf{y})}.$$

where $f(\boldsymbol{\xi}_{-0}|\boldsymbol{\xi}_0) = \int_{\boldsymbol{\pi}} f(\boldsymbol{\xi}_{-0}|\boldsymbol{\xi}_0, \boldsymbol{\pi}) f(\boldsymbol{\pi}) d\boldsymbol{\pi}$ and $f(\boldsymbol{\xi}_0) = \int_{\boldsymbol{\pi}_0} f(\boldsymbol{\xi}_0|\boldsymbol{\pi}_0) f(\boldsymbol{\pi}_0) d\boldsymbol{\pi}_0$ can be computed in closed form: $f(\boldsymbol{\xi}_{-0}|\boldsymbol{\xi}_0) = \frac{\Gamma(K\beta)^K \prod_{k=1}^K \prod_{h=1}^K \Gamma(n_{k,h} + \beta)}{\Gamma(\beta)^{K^2} \prod_{k=1}^K \Gamma(n_k - \xi_{T,k} + K\beta)}$ and $f(\boldsymbol{\xi}_0) = \frac{1}{K}$

3.3 Computational Details

Model parameters are estimated with a MCMC algorithm. More precisely the μ_s , γ_s , σ_y^2 s and ξ are simulated with a Gibbs sampler while the remaining parameters require the introduction of a Metropolis step. The full conditionals of μ s and γ s are normal distributions, those of σ_y^2 s are inverse gamma. The full conditionals for the latent variables $\boldsymbol{\xi}_t$, $t \in \mathcal{T}$ are multinomial and the vector of probabilities depends on the entire vector of $\boldsymbol{\xi}_{-t} = \boldsymbol{\xi} \setminus \{\boldsymbol{\xi}_t\}$. More precisely, let s^- and s^+ be the regimes on time $t-1$ and $t+1$, i.e. $\xi_{t-1s^-} = 1$ and $\xi_{t+1s^+} = 1$ respectively, let $n_k^{-t'} = \sum_{\substack{t=0 \\ t \neq t'}}^T \xi_{tk}$ and $n_{k,h}^{-t'} = \sum_{\substack{t=1 \\ t \neq t', t \neq t'+1}}^T \xi_{t-1k} \xi_{t,h}$. If $t \in \mathcal{T} \setminus \{0, T\}$.

$$f(\boldsymbol{\xi}_t|\mathbf{r}, \mathbf{x}, \mathbf{y}, \boldsymbol{\xi}_{-t}, \boldsymbol{\Psi}) \propto \prod_{k=1}^K \frac{\binom{n_{s^-,k}^{-t} + \beta + a_{s^-,k,s^+}}{\xi_{tk}} \binom{n_{k,s^+}^{-t} + \beta}{1 - \xi_{tk}}}{\binom{n_k^{-t} - \xi_{T,k} + K\beta}{\xi_{tk}}} f(r_t, x_t, y_t|\boldsymbol{\xi}_t, \boldsymbol{\Psi})$$

where a_{s^-,k,s^+} assumes value 1 if $s^- = k = s^+$, 0 otherwise, while

$$f(\boldsymbol{\xi}_T|\mathbf{r}, \mathbf{x}, \mathbf{y}, \boldsymbol{\xi}_{-T}, \boldsymbol{\Psi}) \propto \prod_{k=1}^K \binom{n_{s^-,k}^{-T} + \beta}{\xi_{Tk}} f(r_T, x_T, y_T|\boldsymbol{\xi}_T, \boldsymbol{\Psi}).$$

and

$$f(\boldsymbol{\xi}_0|\mathbf{r}, \mathbf{x}, \mathbf{y}, \boldsymbol{\xi}_{-0}, \boldsymbol{\Psi}) \propto \prod_{k=1}^K \frac{\binom{n_{k,s^+}^{-0} + \beta}{\xi_{0k}}}{\binom{n_k^{-0} - \xi_{T,k} + K\beta}{\xi_{0k}}}.$$

It is well known that the MCMC sampler for HMM tends to mix really slow (Andrieu *et al.*, 2010). To speed up the convergence we try to find an optimal proposal distribution for the Metropolis step, that samples $K \sigma_1^2$ variables, $K \rho$ variables and $T r$ variables, using the algorithm described in Robert and Casella (2009), page 258. With the goal to speed up the MCMC convergence, as a general advice, is suggested to decrease the dimension of the parameters space, i.e. do as much marginalization as possible (Banerjee *et al.*, 2004). In our model we found convenient to marginalize over the vectors $\boldsymbol{\pi}_k$, $k = 1, \dots, K$ and $\boldsymbol{\pi}_0$ but not over \mathbf{r} . Marginalization over \mathbf{r} decreases significantly the number of random variables to simulate but does not allow to have closed form for full conditional distributions of γ_{k1} , γ_{k2} and $\boldsymbol{\mu}_k$. Without employing the Gibbs step, the MCMC algorithm becomes considerably slower in moving toward its stationary distribution and then the computational burden increases as a larger number of iterations is required. On the other hand, marginalization over $\boldsymbol{\pi}_k$, $k = 1, \dots, K$ and $\boldsymbol{\pi}_0$ has impact only on the way we simulate ξ_t , $t = 0, 1, \dots, T$, but their simulation is simple in both cases, with or without $\boldsymbol{\pi}_k$, $k = 1, \dots, K$ and $\boldsymbol{\pi}_0$, and can be carried out in a Gibbs step.

In the estimation step, we take into account the label-switching issue, common to all latent-class-based models. This problem occurs when exchangeable priors are used for the state specific parameters, which is common practice if there are not prior informations about the hidden states. In these cases, the posterior distribution is invariant to permutations of the state labels and, hence, the marginal posterior distributions of the state specific parameters are identical for all states. Therefore, direct inferences about the state specific parameters are not available from the MCMC output. Various approaches to deal with the label switching problem in finite mixture models have been proposed in the literature; see Jasra *et al.* (2005) for a recent review. To tackle the label switching we decide to use the post processing technique called *pivotal reordering*, proposed in Spezia (2009) or in Marin and Robert (2013), Chapter 6.5.

3.4 Model Selection

To decide the number of regimes, we considered the idea of use the Reversible Jump (Green, 1995) or a non parametric approach, as the one proposed by Teh *et al.* (2004). However, our main goal is to demonstrate that the CL-GPN is suitable in a HMM Bayesian framework to model circular-linear variables. Thus, we do not want to further increase the complexity of an already highly complex model by introducing K as random variable.

Common model choice criteria are AIC, BIC, ICL and different classification-based information criteria which are minimized among a set of potential models. We evaluate these criteria using the set of parameters, among the MCMC draws, that maximize the posterior distribution (called maximum a posterior, MAP or MAP estimator) (Frühwirth Schnatter, 2006, Section 4.4.2, 7.1.4). Let $\tilde{\Psi}$ be the MAP estimator, we compute the *BIC* and *AIC* as $BIC = -2 \log \left(f(\mathbf{x}, \mathbf{y} | \tilde{\Psi}) \right) + \#parameters \times \log(T)$ and $AIC = -2 \log \left(f(\mathbf{x}, \mathbf{y} | \tilde{\Psi}) \right) + 2 \times \#parameters$.

The BIC and AIC are generally criticized since they do not take into account the quality of classification of the variables in the K regimes. For classification purpose Biernacki *et al.* (2000) propose to use the ICL; an index based on the likelihood of observed variables and the vector of regimes indicator that is used by Celeux and Durand (2008) in a HMM context. We compute a BIC approximation of the $ICL = f(\mathbf{x}, \mathbf{y} | \tilde{\xi}, \tilde{\Psi}) - 2 \log f(\tilde{\xi}) + \#parameters \times \log(T)$, (see for example Frühwirth Schnatter, 2006, pag. 214)

in the latter case, as suggested by McLachlan and Peel (2000), pag. 216, we first obtain an estimator of ξ , i.e. the MAP $\tilde{\xi}$, and then, as for the BIC and AIC, we compute the ICL using the MAP estimator of Ψ conditioning on the value $\tilde{\xi}$.

4 SIMULATION STUDY

In this Section we carried out a simulation study to investigate the performance of the proposed approach in recovering model parameters and the hidden structure of the data. We empirically demonstrate that the CL-GPN can be used in presence of both unimodal or bimodal state-dependent circular distributions and that ignoring the dependencies between the circular and linear variable at a given time leads to a higher number of states.

We plan the simulation study to cover schemes with different underlying *null* models assuming bimodal or almost uniform shapes for the circular variable, and overlapping or well-separated state-dependent distributions for the linear variable. On each simulated datasets, we estimate three models: 1) the CL-GPN model; 2) a constrained model, defined as diagonal CL-GPN (CL-DPN), with $\Sigma_k = \mathbf{I}_2$, so that the state-dependent circular distribution is symmetric and unimodal; 3) a CL-GPN model with all the γ_{k1} and γ_{k2} equal to zero, i.e. assuming independence between circular and linear variable given the latent state (indicated as Ind-CL-GPN).

4.1 Designing the simulation study

For each null model, we simulated 200 datasets considering two time-series lengths, $T = 500$ and $T = 2000$, with $K = 3$, $\xi_0 = 1$ and transition matrix π with diagonal elements equal to 0.8 and extradiagonal elements equal to 0.1. The considered schemes are summarized in Figure 2, and are characterized by the following settings:

- a) distributions C1 and L1, C2 and L2 and C3 and L3 are considered as state-dependent distributions for the first, the second and the third regime, respectively. The joint representation through scatters is displayed in Figure 3. This scheme has bimodal state-dependent circular distributions and well separated linear ones. The following parameters are used to generate data:

$$\boldsymbol{\mu} = \begin{bmatrix} \mu_{11} & \mu_{12} & \mu_{13} \\ \mu_{21} & \mu_{22} & \mu_{23} \end{bmatrix} = \begin{bmatrix} 0.1 & 0.1 & 0.0 \\ 0.1 & -1.0 & -0.1 \end{bmatrix}$$

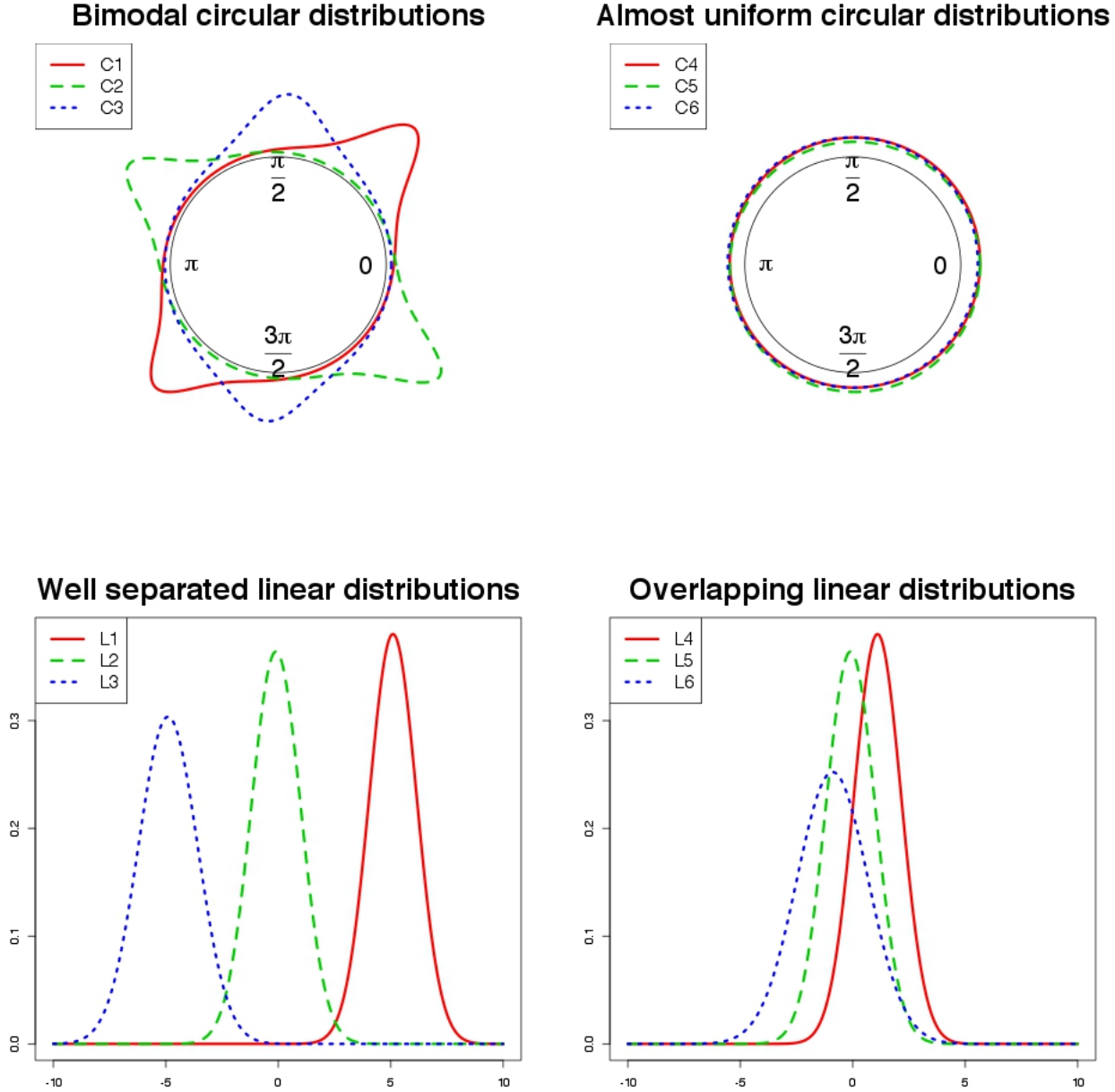


Figure 2: Marginal distributions used in the simulation examples.

$$\gamma = \begin{bmatrix} \gamma_{01} & \gamma_{02} & \gamma_{03} \\ \gamma_{11} & \gamma_{12} & \gamma_{13} \\ \gamma_{21} & \gamma_{22} & \gamma_{23} \end{bmatrix} = \begin{bmatrix} 5.0 & 0.0 & -5.0 \\ 1.0 & 0.0 & 1.0 \\ 0.0 & -1.0 & 1.0 \end{bmatrix}$$

$$\sigma_{k1}^2 = \begin{cases} 1 & k=1 \\ 2 & k=2 \\ 0.1 & k=3 \end{cases} ; \quad \sigma_{ky}^2 = \begin{cases} 0.1 & k=1 \\ 0.2 & k=2 \\ 0.5 & k=3 \end{cases} ; \quad \rho_k = \begin{cases} 0.9 & k=1 \\ -0.9 & k=2 \\ 0.2 & k=3 \end{cases}$$

b) this setting shares circular distributions with scheme a) while the state-dependent linear distributions are respectively the density L4, L5 and L6 of Figure 2. The joint representation through scatters is displayed in Figure 3. With respect to scheme a), we change the values of γ :

$$\gamma = \begin{bmatrix} \gamma_{01} & \gamma_{02} & \gamma_{03} \\ \gamma_{11} & \gamma_{12} & \gamma_{13} \\ \gamma_{21} & \gamma_{22} & \gamma_{23} \end{bmatrix} = \begin{bmatrix} 1.0 & 0.0 & -1.0 \\ 1.0 & 0.0 & 1.0 \\ 0.0 & -1.0 & 1.0 \end{bmatrix}$$

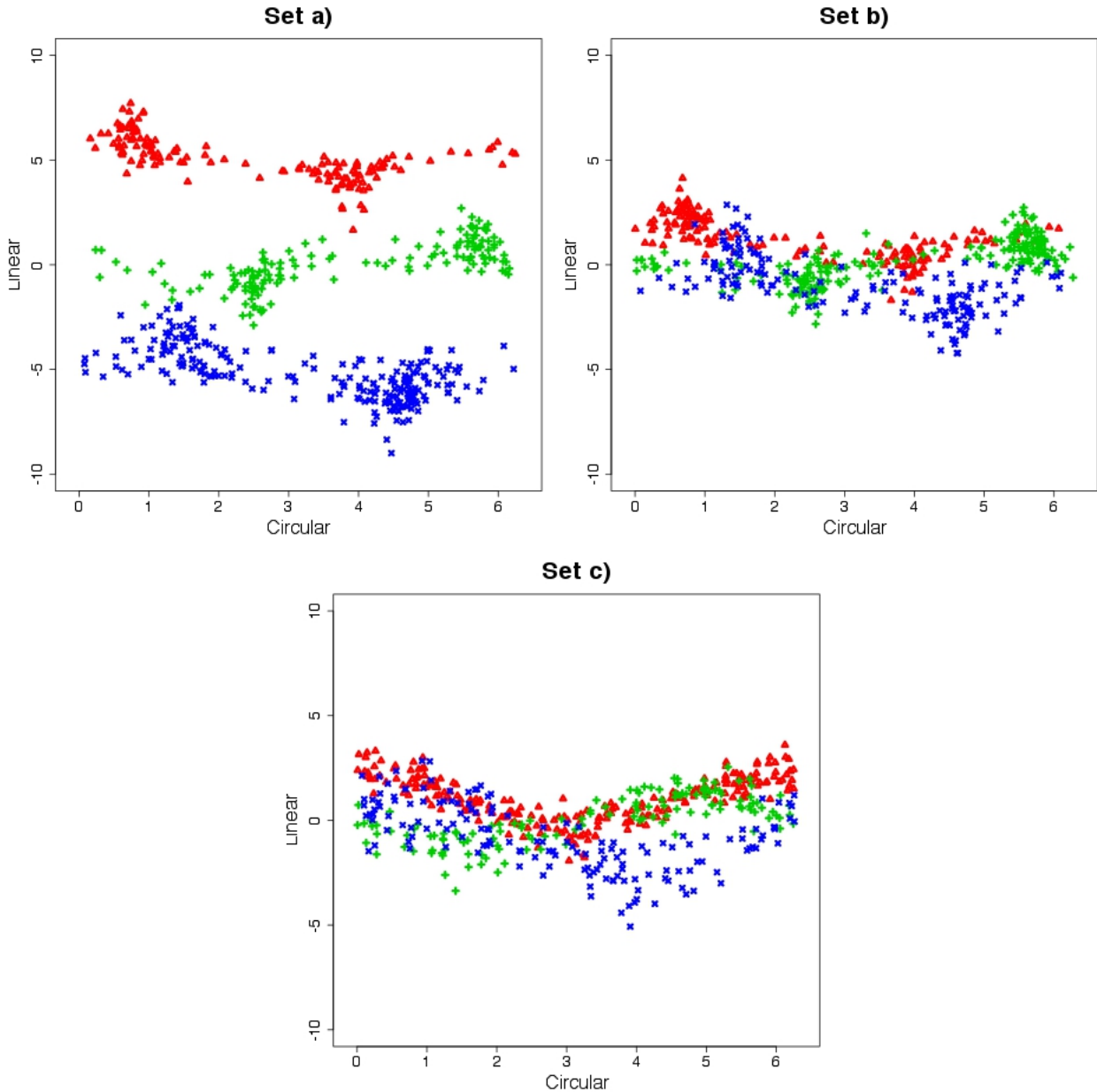


Figure 3: Scatter plot of one simulated dataset for each set of parameters ($T = 500$).

to have more overlapping state-dependent distributions for the linear variable.

- c) the state-dependent distributions for the linear variable are the same as in scheme b), whilst the circular ones are respectively the density C4, C5 and C6 of Figure 2. The joint representation through scatters is displayed in Figure 3. In this case we simulate from a CL-DPN since we use $\sigma_{11}^2 = \sigma_{12}^2 = \sigma_{13}^2 = 1$ and $\rho_1 = \rho_2 = \rho_3 = 0$, i.e. the circular variable has state-dependent unimodal (almost uniform) distributions.

On each dataset we estimate models with K from 2 to 6 and assuming the following prior distributions: $\mu_{ki} \sim N(0, 5)$, $\gamma_{kj} \sim N(0, 5)$, $\rho_k \sim N(0, 5)I(-1, 1)$, $\sigma_{k1}^2 \sim IG(2, 1)$ ¹, $\sigma_{ky}^2 \sim IG(2, 1)$, $\beta = 1$ with $i = 1, 2$ and $j = 1, 2, 3$; i.e. they do not depend on the regime.

¹The two parameters are the shape and rate, respectively

Table 1: Frequency distribution of predicted number of regimes

T	Model	Scheme	Predicted K (AIC)					Predicted K (BIC)					Predicted K (ICL)				
			2	3	4	5	6	2	3	4	5	6	2	3	4	5	6
500	CL-GPN	a)	0.00	0.62	0.35	0.03	0.00	0.00	0.82	0.18	0.00	0.00	0.00	0.91	0.09	0.01	0.00
500	CL-GPN	b)	0.00	0.68	0.25	0.07	0.00	0.00	0.67	0.30	0.03	0.00	0.00	0.90	0.09	0.01	0.00
500	CL-GPN	c)	0.00	0.98	0.02	0.00	0.00	0.00	0.98	0.02	0.00	0.00	0.00	0.99	0.01	0.00	0.00
500	CL-DPN	a)	0.00	0.59	0.37	0.04	0.00	0.00	0.80	0.18	0.02	0.00	0.00	0.87	0.10	0.03	0.00
500	CL-DPN	b)	0.00	0.61	0.31	0.08	0.00	0.00	0.63	0.33	0.04	0.00	0.00	0.86	0.12	0.02	0.00
500	CL-DPN	c)	0.00	0.98	0.01	0.01	0.00	0.00	1.00	0.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00
500	Ind-CL-GPN	a)	0.00	0.00	0.01	0.08	0.91	0.00	0.00	0.00	0.09	0.91	0.00	0.00	0.08	0.26	0.66
500	Ind-CL-GPN	b)	0.00	0.00	0.02	0.08	0.90	0.00	0.00	0.05	0.03	0.92	0.03	0.00	0.00	0.07	0.90
500	Ind-CL-GPN	c)	0.00	0.00	0.00	0.09	0.91	0.00	0.00	0.00	0.09	0.91	0.48	0.41	0.06	0.03	0.02
2000	CL-GPN	a)	0.00	0.97	0.03	0.00	0.00	0.00	0.98	0.01	0.01	0.00	0.00	0.98	0.00	0.00	0.00
2000	CL-GPN	b)	0.00	0.97	0.01	0.02	0.00	0.00	0.98	0.02	0.00	0.00	0.00	0.98	0.02	0.00	0.00
2000	CL-GPN	c)	0.00	0.96	0.03	0.01	0.00	0.00	0.99	0.01	0.00	0.00	0.00	1.00	0.00	0.00	0.00
2000	CL-DPN	a)	0.00	0.90	0.08	0.02	0.00	0.00	0.91	0.05	0.04	0.00	0.00	0.93	0.07	0.00	0.00
2000	CL-DPN	b)	0.00	0.91	0.07	0.02	0.00	0.00	0.93	0.03	0.04	0.00	0.00	0.91	0.09	0.00	0.00
2000	CL-DPN	c)	0.00	0.98	0.02	0.00	0.00	0.00	1.00	0.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00
2000	ind-CL-GPN	a)	0.00	0.00	0.02	0.31	0.67	0.00	0.00	0.02	0.27	0.71	0.00	0.04	0.21	0.31	0.44
2000	ind-CL-GPN	b)	0.00	0.00	0.03	0.07	0.90	0.00	0.00	0.02	0.08	0.90	0.03	0.06	0.22	0.41	0.28
2000	ind-CL-GPN	c)	0.00	0.00	0.00	0.29	0.71	0.00	0.00	0.00	0.25	0.75	0.06	0.27	0.39	0.24	0.04

4.2 Simulation Study Results

To evaluate the performance of AIC, BIC and ICL as selection criteria for the number of regimes, in Table 1 we report the frequency distribution of the predicted K under each simulation setting considered for the CL-GPN, CL-DPN and Ind-CL-GPN models. With respect to the CL-GPN model, we can observe that ICL performs considerably well in all cases. In fact, the predicted K is only occasionally different from the true one, and, when this happens, the former is always larger than the latter. On the other hand, AIC and BIC have an excellent behavior with the exception of the cases $T = 500$, schemes a) and b). As may be expected, these criteria perform better as the amount of information in the data increases.

Ignoring the (state-dependent) correlation between circular and linear measurements may strongly affect the hidden structure. Indeed, by looking at information criteria for the Ind-CL-GPN model, we have that the latent structure is not well recovered and a higher number of regimes than expected is estimated. Of course, this affects parameter estimates and results interpretation, as a not needed number of (latent) regimes is identified in the data.

Here we briefly summarize the results of the simulation study for scheme c). By looking at parameters estimates (see Table 2), we have that the CL-GPN and the CL-DPN models lead essentially to the same results. Point estimates and credibility intervals are very close, suggesting that the CL-GPN distribution can be used whenever we cast doubts on the unimodality of circular distributions. Indeed, the CL-DPN distribution is a specific case of the CL-GPN one, in which conditional circular distribution are constrained to be unimodal.

To further resemble empirical situations, we randomly drop 10% observation of a randomly selected dataset simulated accordingly to the scheme c) with $T = 500$ and estimate the CL-GPN and a CL-DPN models. Along with model parameters, we also simulate the missing observations. We compute the average continuous ranked probability score (CRPS) for both the circular (Grimit *et al.*, 2006) and linear variable (Gneiting and Raftery, 2007) from the posterior samples of the missing observations, as well as the average prediction error (APE) for the circular variable (Jona Lasinio *et al.*, 2012) and the mean squared error for the linear ones (MSE). With the CRPS we evaluate the model performance regarding the entire predictive distribution. APE and MSE allow us to measure the distance between the true values and the simulated ones. The CPRPs for the circular variable and the MSE are identical under the two models while the CRPS for the linear one is 0.66 under the CL-GPN and 0.67 under the CL-DPN and the APE is respectively 0.76 and 0.75 for the CL-GPN and the CL-DPN. Then the two models have the same performances in dealing with the missing values as well.

From the computational point of view, in the datasets with $T = 2000$ our C++ implementation of the model needs 1000000 iterations with a burnin of 700000 and a thin of 100 while with $T = 500$ the iterations needed are 800000, with a burnin of 400000 and again a thin of 100. The computational work for the simulation study and the real data application of Section 5, has been executed on the IT resources made available by ReCaS, a project financed by the MIUR (Italian Ministry for Education, University and Research) in the ‘‘PON Ricerca e Competitivit 2007-2013 - Azione I - Interventi di rafforzamento strutturale’’ PONa3_00052, Avviso 254/Ric. The computational time are of the order

Table 2: posterior median estimates of the parameter ($\hat{\cdot}$) and credibility intervals (CI) for scheme c) and $T = 500$

	CL-GPN			CL-DPN		
	$k = 1$	$k = 2$	$k = 3$	$k = 1$	$k = 2$	$k = 3$
$\hat{\mu}_{k1}$	0.12	0.12	-0.05	0.13	0.11	-0.06
CI	(-0.04 0.28)	(-0.09 0.33)	(-0.26 0.15)	(-0.03 0.31)	(-0.11 0.34)	(-0.27 0.16)
$\hat{\mu}_{k2}$	0.04	-0.08	0.07	0.03	-0.09	0.09
CI	(-0.12 0.19)	(-0.30 0.12)	(-0.14 0.27)	(-0.13 0.18)	(-0.32 0.14)	(-0.12 0.29)
$\hat{\rho}_k$	0.06	-0.09	0.14	.	.	.
CI	(-0.13 0.22)	(-0.32 0.16)	(-0.11 0.37)	(. .)	(. .)	(. .)
$\hat{\sigma}_{k1}^2$	0.94	0.76	1.01	.	.	.
CI	(0.65 1.35)	(0.47 1.22)	(0.63 1.57)	(. .)	(. .)	(. .)
$\hat{\gamma}_{k0}$	0.98	-0.04	-1.04	0.98	-0.04	-1.05
CI	(0.89 1.06)	(-0.19 0.13)	(-1.22 -0.86)	(0.89 1.06)	(-0.20 0.12)	(-1.23 -0.88)
$\hat{\gamma}_{k1}$	1.04	0.14	0.82	1.01	0.15	0.86
CI	(0.88 1.22)	(-0.04 0.35)	(0.60 1.07)	(0.91 1.13)	(-0.01 0.33)	(0.68 1.06)
$\hat{\gamma}_{k2}$	-0.04	-1	1.04	-0.02	-0.96	1.05
CI	(-0.11 0.05)	(-1.18 -0.83)	(0.82 1.28)	(-0.09 0.05)	(-1.12 -0.81)	(0.87 1.26)
$\hat{\sigma}_{ky}^2$	0.14	0.32	0.42	0.14	0.3	0.41
CI	(0.09 0.19)	(0.18 0.52)	(0.26 0.67)	(0.10 0.19)	(0.18 0.49)	(0.25 0.63)

of 1 hour for $T = 500$ and 5 hours for $T = 2000$. All the results shown are from MCMC chains that reach the convergence, checked using the standard tool on the R package coda.

5 REAL DATA EXAMPLE

Finally, we apply the CL-GPN hidden Markov model to a bivariate time series of wind directions and (log-transformed) speeds. Data are recorded on a semi-hourly base from 12/12/2009 to 12/1/2010 in Ancona (Italy) at a bouy located in the Adriatic Sea 30km from the coast (see Figure 4). Data are recorded on $T = 1500$ times and have been previously analysed by Bulla *et al.* (2012).

As often arise in environmental studies, data are not complete. 213 and 210 missing values are recorded for directions and speeds, respectively; 125 profiles are completely missing.

During wintertime, relevant wind events in the Adriatic Sea are typically generated by the south-eastern Sirocco, the north-eastern Bora and the north-western Maestral. Sirocco arises from a warm, dry, tropical air mass that is pulled northwards by low-pressure cells moving eastwards across the Mediterranean Sea. By contrast, Bora episodes occur when a polar high-pressure area sits over the snow-covered mountains of the interior plateau behind the coastal mountain range and a calm low-pressure area lies further south over the warmer Adriatic. Finally, the Maestral is a sea-breeze wind blowing northwesterly when the east Adriatic coast gets warmer than the sea. While Bora and Sirocco episodes are usually associated with high-speed flows, Maestral is in general linked with good meteorological conditions. Hence, the marginal distribution of (log-transformed) wind speed may be interpreted as the result of mixing different wind-speed regimes.

As for the simulation examples, we look at the AIC, BIC and ICL to select the appropriated number of components. The ICL suggest to use $K = 3$ while the AIC and BIC $K = 4$. To help decide between the two number of regimes, we look at their predictive ability, the CRPS_c and APE highlight loss of predictive ability on the circular variable if we choose $K = 4$ (CRPS_c=0.59 and APE=0.75 with $K = 4$ while CRPS_c=0.34 and APE= 0.35 with $K = 3$). For the linear variable, looking at the values of CRPS_l and MSE, there is a small difference between $K = 3$ and $K = 4$ however both CRPS_l and MSE favour $K = 3$ (CRPS_l=0.17 and APE=0.39 with $K = 4$ while CRPS_l=0.16 and APE= 0.34 with $K = 3$). We decide to adopt $K = 3$, that is also the choice of Bulla *et al.* (2012) following their suggestion that three regimes provide well-separated and more interpretable states. The resulting classification is displayed in Figure 5 and all the credibility intervals and point estimates of the parameters are in Table 3. The estimated transition probabilities are displayed in Table 4. As

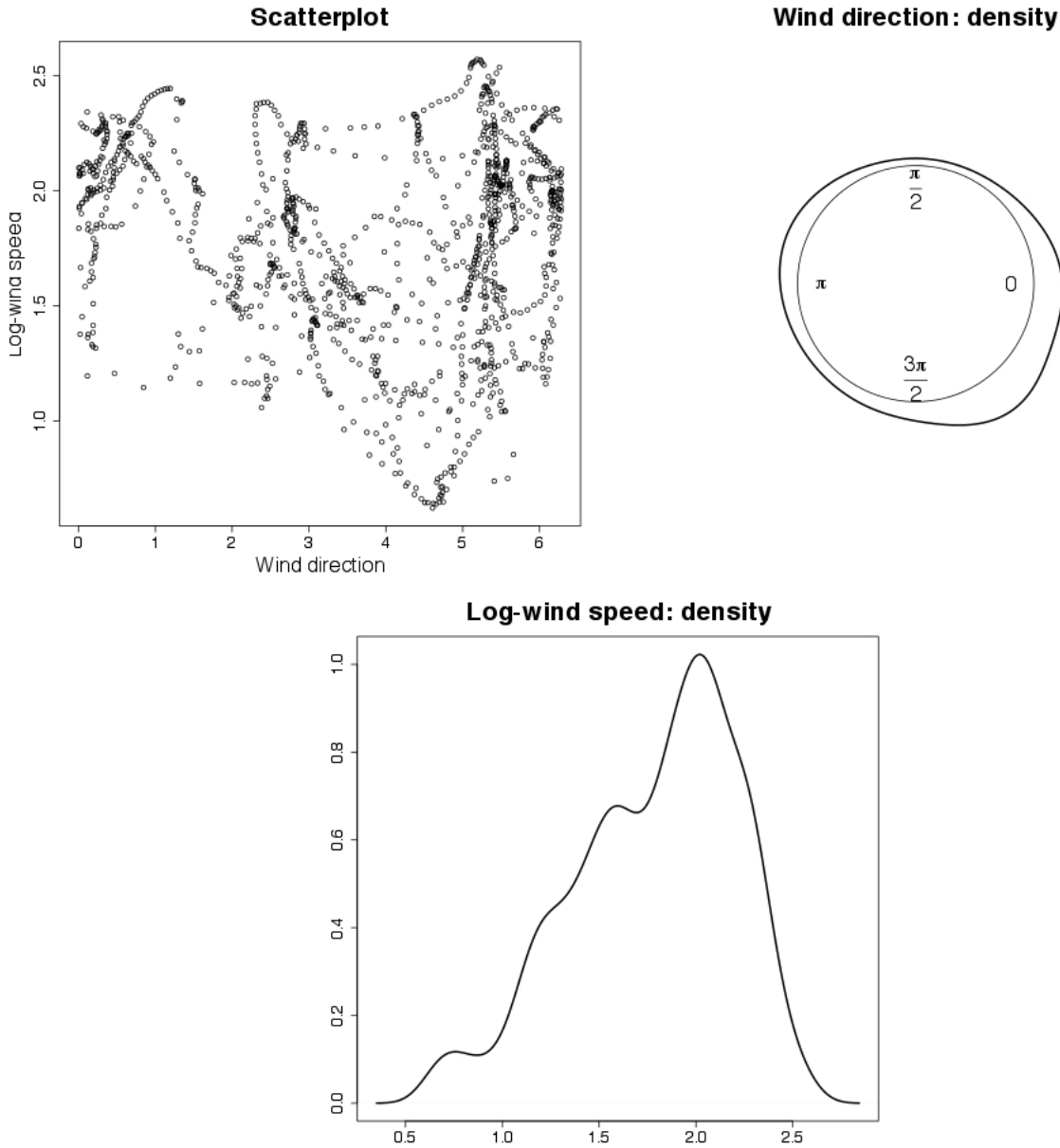


Figure 4: Real data.

expected, the transition probability matrix is essentially diagonal, reflecting the temporal persistence of the regimes, i.e. of wind conditions. Furthermore, the small off-diagonal transition probabilities between states indicate that direct transitions between Sirocco and Bora episodes are very unlikely. The model hence confirms that the Adriatic Sea typically alternates relevant wind events with periods of good conditions.

For a more clear interpretation of the state dependent distributions, we compute some feature of the CL-GPN distribution. In detail, we look at the posterior marginal mean and variance of the linear distribution for each regime ($\hat{\mu}_{ky}$ and $\hat{\sigma}_{ky}^2$), the circular mean ($\hat{\mu}_{kx}$) and concentration (\hat{g}_{kx}) of the circular variable and a measure of correlation between the circular and linear variables ($\hat{\rho}_{kxy}^2$) as in Mardia and Jupp (1999), pag. 245. Point estimates and credibility intervals are provided in Table 5).

The regimes are ordered according to the marginal log-wind speed. In the three regimes the point estimates are respectively $\hat{\mu}_{ky} = 1.42, 1.70, 2.11$ (which correspond to $4.14m/s, 5.47m/s, 8.25m/s$ in the natural scale). With the increases of the velocity, the distribution becomes more concentrated: the marginal linear variance, $\hat{\sigma}_{ky}^2$, is respectively $0.17, 0.11, 0.04$, for a plot of the distributions see

Table 3: Real data: posterior median estimates of the parameter ($\hat{\cdot}$) and credibility intervals (CI)

	$k = 1$	$k = 2$	$k = 3$
$\hat{\mu}_{k1}$	0.45	-1.62	1.19
CI	(0.28 0.63)	(-1.83 -1.41)	(1.02 1.32)
$\hat{\mu}_{k2}$	-1.80	0.67	-0.41
CI	(-2.07 -1.57)	(0.52 0.81)	(-0.51 -0.30)
$\hat{\rho}_k^2$	0.36	0.56	-0.27
CI	(0.14 0.54)	(0.38 0.71)	(-0.46 -0.06)
$\hat{\sigma}_{k1}^2$	2.03	0.94	0.15
CI	(1.46 2.85)	(0.71 1.28)	(0.10 0.23)
$\hat{\gamma}_{k0}$	1.34	1.47	2.36
CI	(1.21 1.49)	(1.33 1.62)	(2.25 2.45)
$\hat{\gamma}_{k1}$	0.01	-0.09	-0.22
CI	(-0.03 0.05)	(-0.16 -0.03)	(-0.30 -0.13)
$\hat{\gamma}_{k2}$	-0.04	0.12	-0.02
CI	(-0.11 0.03)	(0.06 0.18)	(-0.05 0.01)
$\hat{\sigma}_{ky}^2$	0.17	0.09	0.03
CI	(0.14 0.19)	(0.08 0.11)	(0.03 0.06)

Table 4: Real data: Transition matrix

Destination state		1	2	3
Origin state	1	0.96 (0.94, 0.97)	0.02 (0.01, 0.04)	0.02 (0.01, 0.04)
	2	0.02 (0.01 0.04)	0.97 (0.95, 0.98)	0.00 (0.00, 0.01)
	3	0.02 (0.01, 0.03)	0.00 (0.00, 0.01)	0.98 (0.97, 0.99)

Table 5: Real data: posterior median estimates ($\hat{\cdot}$) and credibility intervals (CI) of the features of the distribution CL-GPN

	$k = 1$	$k = 2$	$k = 3$
$\hat{\mu}_{kx}$	4.98	2.7	6.04
CI	(4.81 5.16)	(2.55 2.83)	(5.90 6.19)
\hat{g}_{kx}	0.78	0.83	0.82
CI	(0.71 0.84)	(0.77 0.87)	(0.76 0.86)
$\hat{\mu}_{ky}$	1.42	1.70	2.11
CI	(1.34 1.50)	(1.63 1.78)	(2.02 2.20)
$\hat{\sigma}_{ky}^2$	0.17	0.11	0.04
CI	(0.15 0.20)	(0.09 0.12)	(0.03 0.07)
$\hat{\rho}_{kxy}^2$	0.02	0.05	0.05
CI	(0.00 0.10)	(0.00 0.18)	(0.00 0.17)

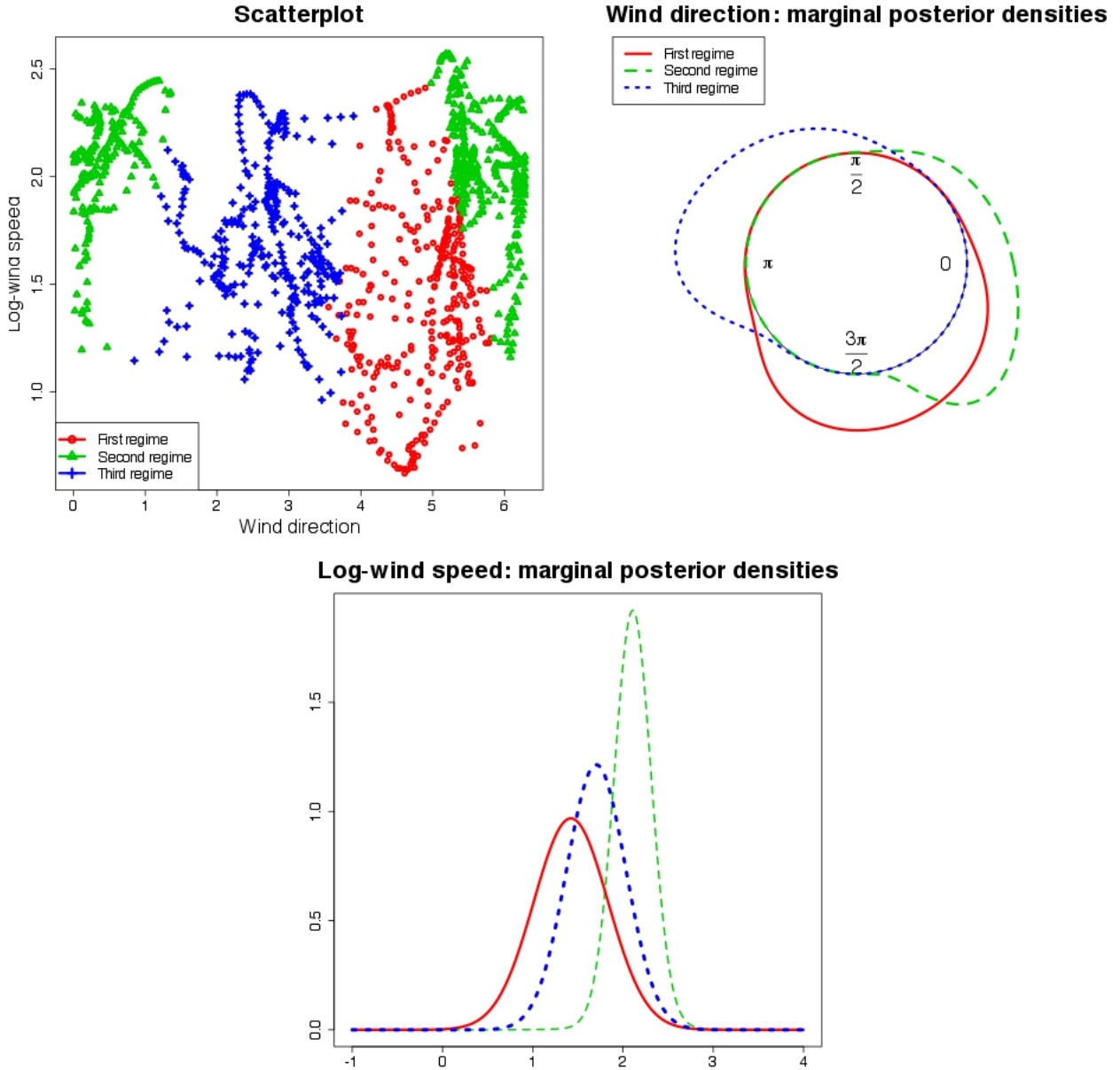


Figure 5: Real data classification.

Figure 5. The circular mean is 4.98 in the first regime, north-westerly Maestral episodes, 2.70 in the second, south-eastern Sirocco, and 6.04 in the third, northern Bora jets, on the first regime the circular marginal distribution is less concentrated than in the others (0.78 for $k = 1$, 0.83 for $k = 2$ and 0.82 for $k = 3$).

The correlations between the circular and linear variables are weak in all the regimes: $\hat{\rho}_{kxy}^2$ is 0.02 in the first and 0.05 in the others. Under the hypothesis of no correlation, i.e. $\hat{\rho}_{kxy}^2 = 0$, the statistic $\tilde{F} = \frac{\hat{\rho}_{kxy}^2(n-1)}{1-\hat{\rho}_{kxy}^2}$ is distributed as a $F_{2,T-3}$ where in our case $T = 1500$. The lower limits of the 95% credibility intervals of the posterior distributions of \tilde{F} are 1.24, 4.80 and 4.95 in the calm, transition and storm conditions respectively, and the 95% percentile of $F_{2,T-3}$ is 3.00. Accordingly, circular-linear correlations are significant in the transition and storm conditions only. This result is not present at all in previous analyses. This can be seen also with the value of γ_{k1} and γ_{k2} in Table 3. In the first regime $\hat{\gamma}_{11} = 0.01$ and $\hat{\gamma}_{12} = -0.04$, both credibility intervals contain the 0. In the second there is a negative relation between the linear variable and the cosine of the circular one ($\gamma_1 = -0.09$) and a positive relation with the sine ($\gamma_2 = 0.12$). In the third regime the dependence between the

linear and circular variable is on the cosine direction ($\gamma_1 = -0.22$).

We estimate the model using 1000000 iterations, a burnin of 700000 and a thin of 100. Here again we checked the convergence of the MCMC chain using the standard tool on the R package coda.

6 Discussion

In this work we introduce, for the first time, the CL-GPN distribution in a Bayesian HMM framework and we present the explicit expression of the CL-GPN likelihood. This approach allows to easily model multivariate processes with mixed support (circular-linear), by combining the bivariate representation of the circular component (i.e. the projected normal distribution) and a Gaussian distribution for the linear part. Here we considered one circular and one linear variable although it is fairly easy to extend the proposed model to more than one linear component.

The Bayesian framework allows us to overcome identifiability issues and computational problems that may arise in the classical setting. Several implementation novelties are introduced to speed up algorithms convergence. We use an adaptive Metropolis whenever a Gibbs sampler is not implementable (section 3.3). Furthermore we marginalize the transition matrix so to avoid its estimation to reduce the problem size obtaining it as an a posteriori byproduct (Section 3.2) and we provide evidence that the marginalization does not affect parameters estimation. We also demonstrate that assuming conditional independence between the circular and linear variable can make difficult to correctly estimate the number of regimes.

We applied this methodology to wind data confirming previously obtained results and highlighting new data features. Circular parameters interpretation is not straightforward, however this does not limit the inferential richness of the model. Using MCMC simulations posterior circular mean and concentration can be derived, as well as the circular-linear correlation. Of course, different areas of application can be considered for the proposed approach, e.g. animal movement modelling (Langrock *et al.*, 2014) and driving behaviour (Jackson *et al.*, 2014).

Further developments will include the extension to more than one circular variable. This extension requires a careful definition of correlation between circular variables that is not straightforward under the projected normal distribution. Another interesting extension of the proposed approach is to allow the estimation of the number of states along with the model parameters. The latter can be obtained using a hierarchical Dirichlet process on the states or a reversible jump.

A crucial assumption of our model is that the temporal dependence is well described by a first order Markov chain, i.e. the sojourn time is geometrical. If we want to allow for different sojourn time distributions with finite support the HMM formulation is exact. Similarly by allowing the number of hidden states to grow with the sample size, we can allow for continuous time, i.e. the hidden distribution can be approximated with arbitrary accuracy using the proposed model. This can be seen as a possible solution to computational issues arising with continuous-valued latent models (Langrock *et al.*, 2012b).

References

- Alfò M, Maruotti A, 2010. A hierarchical model for time dependent multivariate longitudinal data. In Palumbo F, Lauro CN, Greenacre MJ (eds.), *Data Analysis and Classification*, Studies in Classification, Data Analysis, and Knowledge Organization, Springer Berlin Heidelberg, 271–279.
- Andrieu C, Doucet A, Holenstein R, 2010. Particle Markov chain Monte Carlo methods. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* **72**(3): 269–342.
- Banerjee S, Gelfand AE, Carlin BP, 2004. *Hierarchical Modeling and Analysis for Spatial Data*. Chapman and Hall/CRC.
- Bartolucci F, Farcomeni A, 2009. A multivariate extension of the dynamic logit model for longitudinal data based on a latent Markov heterogeneity structure. *Journal of the American Statistical Association* **104**(486): 816–831.

- Bartolucci F, Farcomeni A, 2010. A note on the mixture transition distribution and hidden Markov models. *Journal of Time Series Analysis* **31**(2): 132–138.
- Bartolucci F, Farcomeni A, Pennoni F, 2012. *Latent Markov Models for Longitudinal Data*. Chapman and Hall.
- Bartolucci F, Farcomeni A, Pennoni F, 2014. Latent Markov models: a review of a general framework for the analysis of longitudinal data with covariates. *TEST to appear*.
- Bartolucci F, Pennoni F, Vittadini G, 2011. Assessment of school performance through a multilevel latent Markov Rasch model. *Journal of educational and behaviour statistics* **36**(4): 491–522.
- Baudry JP, Raftery AE, Celeux G, Lo K, Gottardo R, 2010. Combining mixture components for clustering. *Journal of Computational and Graphical Statistics* **19**(2): 332–353.
- Biernacki C, Celeux G, Govaert G, 2000. Assessing a mixture model for clustering with the integrated completed likelihood. *IEEE Trans. Pattern Anal. Mach. Intell.* **22**(7): 719–725.
- Bulla J, Lagona F, Maruotti A, Picone M, 2012. A multivariate hidden Markov model for the identification of sea regimes from Incomplete skewed and circular time series. *Journal of Agricultural, Biological, and Environmental Statistics* **17**(4): 544–567.
- Cappé O, Moulines E, Ryden T, 2005. *Inference in Hidden Markov Models*. Springer Series in Statistics, Springer.
- Celeux G, Durand JB, 2008. Selecting hidden Markov model state number with cross-validated likelihood. *Comput. Stat.* **23**(4): 541–564.
- Dannemann J, 2012. Semiparametric hidden Markov models. *Journal of Computational and Graphical Statistics* **21**(3): 677–692.
- Frühwirth Schnatter S, 2006. *Finite Mixture and Markov Switching Models*. Springer Verlag.
- Geweke J, Amisano G, 2011. Hierarchical Markov normal mixture models with applications to financial asset returns. *Journal of Applied Econometrics* **26**(1): 1–29.
- Gneiting T, Raftery AE, 2007. Strictly proper scoring rules, prediction, and estimation. *Journal of the American Statistical Association* **102**(477): 359–378.
- Green PJ, 1995. Reversible jump Markov chain monte carlo computation and Bayesian model determination. *Biometrika* **82**(4): 711–732.
- Grimit EP, Gneiting T, Berrocal VJ, Johnson NA, 2006. The continuous ranked probability score for circular variables and its application to mesoscale forecast ensemble verification. *Quarterly Journal of the Royal Meteorological Society* **132**: 2925–2942.
- Holzmann H, Munk A, Suster M, Zucchini W, 2006. Hidden Markov models for circular and linear-circular time series. *Environmental and Ecological Statistics* **13**(3): 325–347.
- Jackson J, Albert P, Zhiwei Z, 2014. A two-state mixed hidden Markov model for risky teenage driving behavior. *The Annals of Applied Statistics* **To appear**.
- Jammalamadaka SR, SenGupta A, 2001. *Topics in Circular Statistics*. World Scientific.
- Jasra A, Holmes CC, Stephens DA, 2005. Markov chain Monte Carlo methods and the label switching problem in Bayesian mixture modeling. *Statistical Science* **20**(1): 50–67.
- Johnson RA, Wehrly TE, 1978. Some Angular-Linear Distributions and Related Regression Models. *Journal of the American Statistical Association* **73**(363): 602–606.

- Jona Lasinio G, Gelfand A, Jona Lasinio M, 2012. Spatial analysis of wave direction data using wrapped Gaussian processes. *Annals of Applied Statistics* **6**(4): 1478–1498.
- Kato S, Shimizu K, Shieh GS, 2008. A circular-circular regression model. *Statistica Sinica* **18**: 633–645.
- Lagona F, Maruotti A, Picone M, 2011. *A Non-homogeneous Hidden Markov Model For The Analysis Of Multi-pollutant Exceedances Data — InTechOpen*, chapter 10. INTECH, University Campus STeP Ri Slavka Krautzeka 83/A 51000 Rijeka Croatia, 207–222.
- Lagona F, Picone M, 2011. A latent-class model for clustering incomplete linear and circular data in marine studies. *Journal of Data Science* **9**(4).
- Lagona F, Picone M, 2012. Model-based clustering of multivariate skew data with circular components and missing values. *Journal of Applied Statistics* **39**(5): 927–945.
- Lagona F, Picone M, Maruotti A, Cosoli S, 2014. A hidden Markov approach to the analysis of space-time environmental data with linear and circular components. *Stochastic Environmental Research and Risk Assessment* **to appear**.
- Langrock R, King R, Matthiopoulos J, Thomas L, Fortin D, Morales JM, 2012a. Flexible and practical modeling of animal telemetry data: hidden Markov models and extensions. *Ecology* **93**(11): 2336–2342.
- Langrock R, Kneib T, Michelot T, 2014. Markov-switching generalized additive models. *ArXiv e-prints*.
- Langrock R, MacDonald IL, Zucchini W, 2012b. Some nonstandard stochastic volatility models and their estimation using structured hidden Markov models. *Journal of Empirical Finance* **19**(1): 147–161.
- Langrock R, Swihart BJ, Caffo BS, Punjabi NM, Crainiceanu CM, 2013. Combining hidden Markov models for comparing the dynamics of multiple sleep electroencephalograms. *Statistics in Medicine* **32**(19): 3342–3356.
- Mardia KV, 1976. Linear-Circular Correlation Coefficients and Rhythmometry. *Biometrika* **63**(2).
- Mardia KV, Jupp PE, 1999. *Directional Statistics*. John Wiley and Sons.
- Marin J, Robert C, 2013. *Bayesian Essentials with R*. Springer Texts in Statistics, Springer New York.
- Martinez-Zarzoso I, Maruotti A, 2013. The environmental kuznets curve: functional form, time-varying heterogeneity and outliers in a panel setting. *Environmetrics* **24**(7): 461–475.
- Maruotti A, 2011. Mixed hidden Markov models for longitudinal data: An overview. *International Statistical Review* **79**(3): 427–454.
- McLachlan C, Peel D, 2000. *Finite Mixture Models*. John Wiley and Sons, New York.
- Robert CP, Casella G, 2009. *Introducing Monte Carlo Methods with R*. Springer, Berlin, Heidelberg.
- Rydén T, 2008. Em versus Markov chain Monte Carlo for estimation of hidden Markov models: a computational perspective. *Bayesian Analysis* **3**(4): 659–688.
- Rydén T, Titterton DM, 1998. Computational Bayesian analysis of hidden Markov models. *Journal of Computational and Graphical Statistics* **7**(2): 194–211.
- Spezia L, 2009. Reversible jump and the label switching problem in hidden Markov models. *Journal of Statistical Planning and Inference* **139**(7): 2305 – 2315.
- Spezia L, 2010. Bayesian analysis of multivariate Gaussian hidden Markov models with an unknown number of regimes. *Journal of Time Series Analysis* **31**(1): 1–11.

- Teh YW, Jordan MI, Beal MJ, Blei DM, 2004. Hierarchical Dirichlet processes. *Journal of the American Statistical Association* **101**: 1566–1581.
- Wang F, Gelfand A, 2012. Directional data analysis under the general projected normal distribution. *Statistical Methodology* **10**(1): 113 – 127.
- Wang F, Gelfand A, 2014. Modeling space and space-time directional data using projected Gaussian processes. *Journal of the American Statistical Association* **to appear**.
- Wang F, Gelfand A, Jona Lasinio G, 2014. Joint spatio-temporal analysis of a linear and a directional variable: Space-time modeling of wave heights and wave directions in the adriatic sea. *Statistica Sinica* **to appear**.
- Yildirim S, Singh SS, Dean T, Jasra A, 2014. Parameter estimation in hidden Markov models with intractable likelihoods using sequential monte carlo. *Journal of Computational and Graphical Statistics* **to appear**.
- Zhang Q, Snow Jones A, Rijmen F, Ip EH, 2010. Multivariate discrete hidden Markov models for domain-based measurements and assessment of risk factors in child development. *Journal of Computational and Graphical Statistics* **19**(3): 746–765.
- Zucchini W, MacDonald I, 2009. *Hidden Markov Models for Time Series: An Introduction Using R*. Chapman & Hall/CRC Monographs on Statistics & Applied Probability, Taylor & Francis.