

Clustering of bivariate satellite time series: a quantile approach

Victor Muthama Musau *

Department of Pure and Applied Sciences,
Kirinyaga University, Kenya

Carlo Gaetan

Dipartimento di Scienze Ambientali, Informatica e Statistica,
Università Ca' Foscari - Venezia, Italy

Paolo Girardi

Dipartimento di Scienze Ambientali, Informatica e Statistica,
Università Ca' Foscari - Venezia, Italy

July 26, 2022

Abstract

Clustering has received much attention in Statistics and Machine learning with the aim of developing statistical models and autonomous algorithms which are capable of acquiring information from raw data in order to perform exploratory analysis. Several techniques have been developed to cluster sampled univariate vectors only considering the average value over the whole period and as such they have not been able to explore fully the underlying distribution as well as other features of the data,

*Address for correspondence: Victor Muthama Musau, Department of Pure and Applied Sciences Kirinyaga University, KENYA. E-mail: vmusau@kyu.ac.ke.

especially in presence of structured time series. We propose a model-based clustering technique that is based on quantile regression permitting us to cluster bivariate time series at different quantile levels. We model the within cluster density using asymmetric Laplace distribution allowing us to take into account asymmetry in the distribution of the data. We evaluate the performance of the proposed technique through a simulation study. The method is then applied to cluster time series observed from Glob-colour satellite data related to trophic status indices with aim of evaluating their temporal dynamics in order to identify homogeneous areas, in terms of trophic status, in the Gulf of Gabes.

Keywords: Asymmetric Laplace distribution, model-based clustering, quantile regression, trophic status, water quality

1 Introduction

The European Community with the European Union Water Framework Directive 2000/60/EC (WFD) indicated a series of trophic status indicators with the scope to monitor the status of the sea-water in order to restore and protect water-bodies from further degradation (Directive *et al.*, 2000; Alikas *et al.*, 2015). The light diffuse attenuation coefficient at 490 nm (KD-490) is an ecologically important water property that provides information about the availability of light to underwater communities which influences ecological processes and biogeochemical cycles in natural waters (Yang *et al.*, 2020). Together with KD-490, the Chlorophyll type-a (Chl-a) is commonly used as a proxy for phytoplankton biomass and as an indicator for eutrophication; high concentration levels may lead to hypoxic or anoxic events while low levels may result in improvement in water quality (Dabuleviciene *et al.*, 2020). The joint use of two complementary indices as the KD-490 and the Chl-a allows to take into account different aspects of water quality ranging from availability of light to underwater communities to the trophic status.

In this context, the classification of areas with a different level of impact may help Institutions to define a program of conservation and environmental protection. The European WFD 2000/60/EC set a series of rules to classify areas at different level of “impact” considering summary statistics (average, geometric mean, or percentile) of a single indicator over a predetermined temporal window (1 year, 5 years, etc, ...) with respect to a reference condition represented by no or very low human pressure (Poikāne *et al.*, 2010). The definition of “reference condition” may be particularly complex and limited to each specific context (Pardo *et al.*, 2012).

Measures of the KD-490 and Chl-a indices are obtained from satellite sensors. The use of satellite data needs to consider aspects relative to unreliability for different reasons and mainly due to the presence of cloud contamination or malfunctioning of the sensor onboard the satellite; pre-processing procedures are often required to remove some site and measurements before using the data for the application (Alvera-Azcárate *et al.*, 2012; Stafoggia *et al.*, 2017).

In this paper, we concentrate our efforts to overcome the classification based on simple summary statistics considering the temporal component in order to provide more reliable results since it takes into account the time dynamics of a water bodies indicator (i.e. seasonal behaviour, inter-annual variability, etc. . .). In fact, standard clustering techniques were often applied considering the summary statistics of the determinants of interest, and so, potentially valuable information about the temporal behaviours (e.g. trend, peaks, and seasonal patterns) is lost.

In the environmental field time series clustering has gained popularity for grouping time series with similar temporal pattern covering a wide series of applications and approaches (Cazelles *et al.*, 2008; Giraldo *et al.*, 2012; Shi *et al.*, 2013; Finazzi *et al.*, 2015; Haggarty *et al.*, 2015; Gaetan *et al.*, 2017).

Moreover very often the summary statistics suffer from the lack of robustness in presence of contaminated data. In this respect quantile regression appears attractive due to the possibility to overcome these pre-processing issues with the possibility of choosing a particular quantile of interest (Barbosa *et al.*, 2011; Monteiro *et al.*, 2012). Another juncture concerns that most of the published work considered the clustering of univariate response variable (Barbosa *et al.*, 2011; Monteiro *et al.*, 2012) or even a different weight between average value and temporal trend (Li *et al.*, 2016) while some other attempts were performed taking into account the joint distribution of two or more variables of interest or the joint modelling of more quantiles. (Gaetan *et al.*, 2016; Zhang *et al.*, 2019; Sottile and Adelfio, 2019).

Therefore we propose a Bayesian clustering technique to define groups of temporal patterns that are similar considering the quantile of interest on bivariate time series. More in detail, the application regards monthly time series of Chl-a concentrations and KD-490 levels, obtained by satellite data sensors over the Gulf of Gabes (Tunisia), a Mediterranean zone with important biological resources and rich coastal, marine, and freshwater ecosystems. Since the last few decades, due to fast and uncontrolled urbanization and industrialization, the Gulf of Gabes is experiencing an irreversible degradation of the

local coastal area. (Ayadi *et al.*, 2015; El Zrelli *et al.*, 2017).

The paper is organized as follows. The next section illustrates our data on Chl-a and KD490 on the Gulf of the Gabes. In Section 3 we introduce the Bayesian quantile regression with an extension to the bivariate case and proposing our clustering procedure. In Section 4 we present a simulation study for illustrating the performances and peculiarities of the procedure. Section 5 reports the results of the application to the Chlorophyll-a and KD-490 satellite data on the Gulf of the Gabes. In the last section we discuss the relative strengths and weaknesses of our proposal.

2 Chl-a and KD-490 levels in the Gulf of Gabes

The Gulf of Gabes is a Mediterranean zone with important biological resources and rich coastal, marine and freshwater ecosystems. In the last decades the Gulf of Gabes reported several environmental problems mainly due to the presence of human activities associated with over-fishing and seabed trawling while the presence of chemical factories in the Sfax site resulting on a wide wastewater pollution (Aloulou *et al.*, 2012; Rabaoui *et al.*, 2013; Zaghden *et al.*, 2014; Fourati *et al.*, 2018). The results had led to several issues as the local appearance of red tides (Hamza and El Abed, 1994) as well as the changes or decline of the distribution of some marine species (El Kateb *et al.*, 2016; El Zrelli *et al.*, 2018).

The European Community indicates the diffuse Chlorophyll type-a (Chl-a) concentrations as a trophic status indicator of the sea-water. High Chl-a levels may lead to hypoxic or anoxic events. KD-490 indicates how the solar light can penetrate to deeper water and can be used to evaluate the potential disturbance to the water ecosystem.

In this work, we considered two datasets comprising the Chl-a concentrations and the KD-490 levels made available by ACRI (hermes.acri.fr) in the framework of the GlobColour Project (www.globcolour.info). The datasets were formed by monthly values of Chl-a concentrations and KD-490 levels in the Gulf of Gabes from January, 2003 to December, 2011 for a total of 108 time-points. Data were obtained by calibrating

Ocean Colour data provided by different satellite missions, such as MERIS, SeaWiFS and MODIS. For each month, gridded data with 1.8-km resolution are available and a grid of 4,033 points covers the entire Gulf of Gabes. As reported in Figure 1, both Chl-a concentrations and KD-490 levels exhibit a seasonal cyclical pattern, more evident in the case of chlorophyll, with the presence of a peak during the spring period. While the KD-490 is related to the sea water turbidity more or less correlated to environmental factors (heavy rain, wind direction, etc. . .), the Chl-a index reflects the seasonal bloom in vegetation, directly connected to the seasonal variation of the sunlight window and the sea water temperature.

In Figure 2 we report the average levels of Chl-a concentration and KD-490 index for all the sites over the spatial domain.

The highest average levels of Chl-a and KD-490 were reported near the coastal area, in the north-eastern zone (Sfax) and in the southern part (Djerba) of the Gulf. The spatial distribution of the average values between Chl-a and KD-490 appears to have a similar spatial and temporal behaviour, but not everywhere. In fact as shown in Figure 3 the example site 1 exhibits a strong seasonal pattern both in the Chl-a and KD-490 indicators, with a strong correlation between the two time series (Pearson ρ : 66.1%); in the second example site 2 there is a clear difference between the trend of the Chl-a concentration with respect to the KD-490 one as attested by a weak correlation (Pearson ρ : 30.7%).

The dataset is affected by outliers due to measurement errors (Chl-a concentration lower than 0.05 mg/m^{-3} : 1.0%; KD-490 levels lower than 0.05 m^{-1} : 2.8%) and by a different amount of missing observations (Chl-a: 19.1%; KD-490: 1.8%). The presented dataset presents several features which makes the clustering task challenging: a strong seasonality and a high variability, the presence of outliers and a non bell-shaped distribution. The proposed classification may help to assess the trophic status of this area combining the information provided by different temporal Chl-a and KD-490 levels.

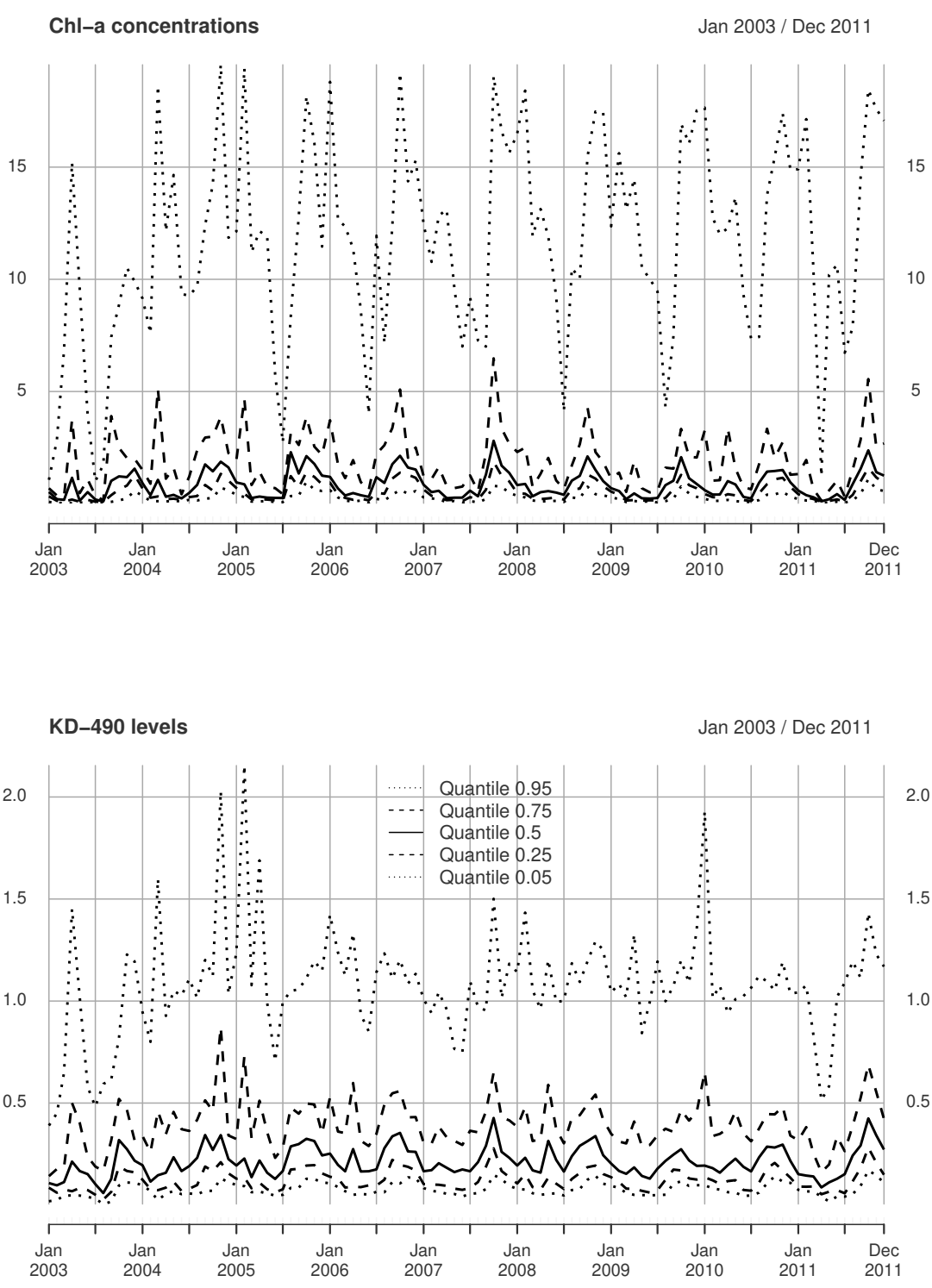


Figure 1: Temporal profile for the quantiles (0.05, 0.25, 0.50, 0.75, and 0.95). For each time we calculated the empirical quantiles over the 4033 time-series.

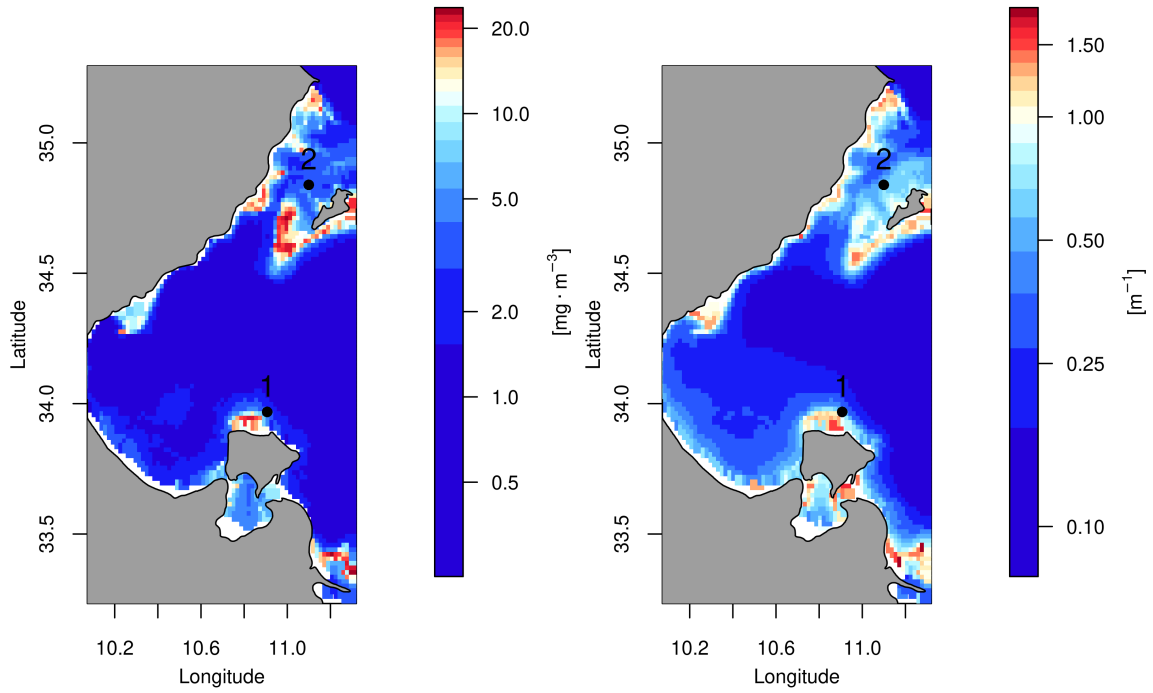


Figure 2: Observed grid-points (4033 sites) in the Gulf of Gabes for the average Chl-a (left) and KD-490 (right) and two example sites reported in Figure 3. White color near the coastal area corresponds to sites with no values due to shallow waters.

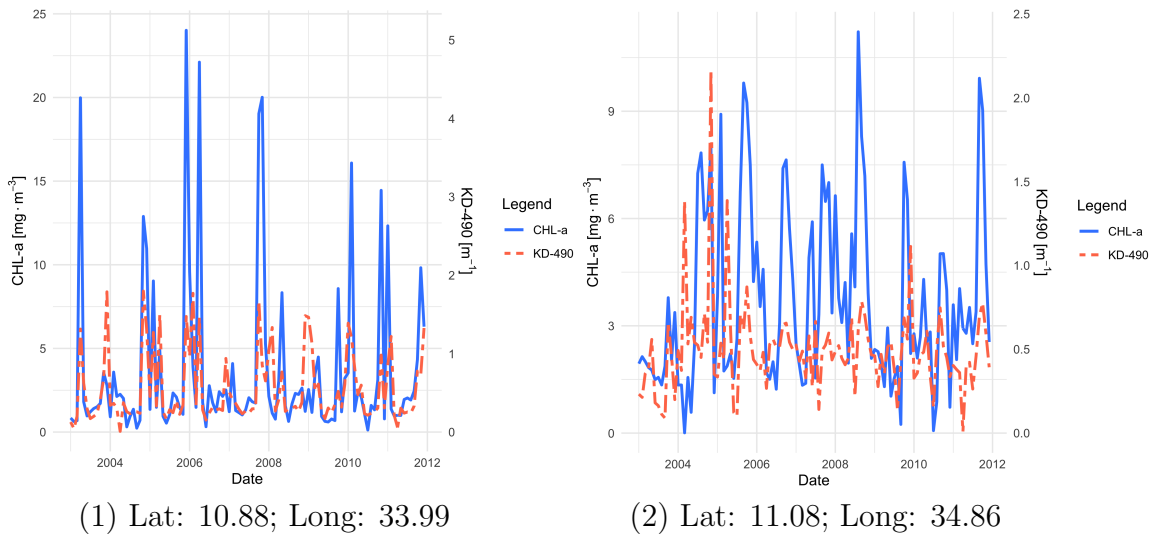


Figure 3: Temporal trend of Chl-a and Kd-490 in the two selected example sites.

3 Statistical modelling and inference

3.1 Multivariate quantile regression and asymmetric Laplace distribution

We start by considering the univariate case. Let $Q_p(y|x)$, for $0 < p < 1$, be the p -th quantile regression function of the univariate continuous random variable y given x , a vector of covariates. We suppose that $Q_p(y|x) = x'\beta$, where β is a vector of unknown parameters to be estimated. Then a quantile regression model can be defined as $y = x'\beta + e$ where e is an error term with density function $f(\cdot; p)$ and the p -th quantile equal to zero, i.e. $\int_{-\infty}^0 f(e; p)de = p$. Owing to the data (y_t, x'_t) , $t = 1, \dots, T$, the estimate of β is classically (Koenker and Bassett, 1978) obtained by minimizing

$$\hat{\beta} = \underset{\beta}{\operatorname{argmin}} \sum_{t=1}^T \rho_p(y_t - x'_t\beta), \quad (1)$$

where $\rho_p(\cdot)$ is the check loss function, i.e. $\rho_p(z) = z(p - I(z < 0))$.

Koenker and Machado (1999) showed that there is a direct relationship between minimizing (1) and the maximum likelihood theory using independent variables y_t with Asymmetric Laplace (AL) density

$$f(e; p) = \frac{p(1-p)}{\sigma} \exp \left\{ -\rho_p \left(\frac{e}{\sigma} \right) \right\}, \quad (2)$$

where $\sigma > 0$ is an additional scale parameter.

In the following we will exploit the representation of y with density (2) as a location scale mixture of Gaussian random variable (Kotz *et al.*, 2001), namely

$$y = x'\beta + \theta\sigma w + \omega\sigma\sqrt{w}\nu \quad (3)$$

where $\nu \sim \mathcal{N}(0, 1)$, and w is a exponential random variable with $E(w) = 1$. Here ν and w are mutually independent and $\theta = (1 - 2p)/\{p(1 - p)\}$ and $\omega^2 = 2/\{p(1 - p)\}$.

The literature has focused on univariate response variable with only a few studies considering extension to multivariate case (Benoit and Van den Poel, 2012; Benoit *et al.*, 2013; Waldmann and Kneib, 2015; Petrella and Raponi, 2019).

In particular Petrella and Raponi (2019) considered a multivariate asymmetric Laplace distribution (Kotz *et al.*, 2001) to specify a quantile regression model for the random vector $\tilde{y} = (y_1, \dots, y_q)'$ in the mixture representation

$$\tilde{y} = X'\tilde{\beta} + D\Theta w + \sqrt{w}D\Sigma^{1/2}\tilde{\nu} \quad (4)$$

Here $\tilde{\nu}$ denotes a q -dimension standard Gaussian vector and w is a exponential random variable with unit mean. The matrix X is a $q \times L$ regressor matrix and $\tilde{\beta}$ is a L -dimensional unknown vector. The other parameters of the model are contained in $D = \text{diag}(\tilde{\sigma})$, with $\tilde{\sigma} = (\sigma_1, \dots, \sigma_q)'$, $\sigma_j > 0$, $j = 1, \dots, q$ and $\Theta = \text{diag}(\theta_1, \dots, \theta_q)'$ with $\theta_j = (1 - 2p_j)/[p_j(1 - p_j)]$.

Moreover the matrix $\Sigma = \Omega R(\tilde{\phi})\Omega$ is a $q \times q$ positive definite matrix with $R(\tilde{\phi})$ being a correlation matrix, that depends on a vector of parameters $\tilde{\phi}$, and $\Omega = \text{diag}(\omega_1^2, \dots, \omega_q^2)$, with entry $\omega_j^2 = 2/[p_j(1 - p_j)]$. Note that using this formulation the component y_j will always be mutually dependent, even though $R(\phi)$ is an identity matrix.

In this paper we consider a slight different approach with respect to (4). We use the same idea as in Waldmann and Kneib (2015) and we set

$$\tilde{y} = X'\beta + D\Theta\tilde{w} + W^{1/2}D\Sigma^{1/2}\tilde{\nu} \quad (5)$$

with $\tilde{w} = (w_1, \dots, w_q)'$ and $W = \text{diag}(w_1, \dots, w_q)$ allowing a different exponential random variable with unit mean w_j for each component.

With this variation model (5) does not define a proper multivariate asymmetric Laplace distribution as in Kotz *et al.* (2001). However the model becomes more flexible and can cover the case of mutually independent components. Now the challenge lies in finding a multivariate density for \tilde{w} , $e(\tilde{w}; \gamma)$, with marginal exponential distributions for the w_j in order to guarantee that the marginal distributions for the response variables are still

asymmetric Laplace. Here the parameter γ is a generic dependence parameter.

Several multivariate exponential distributions have been proposed in the literature (Kotz *et al.*, 2000, Ch. 47). Waldmann and Kneib (2015) adopt the simplest specification that the exponential variables are independent, i.e. $\gamma = 0$.

In view of the bivariate real data example in this paper we exemplify our construction by using the bivariate density proposed by Downton (1970), namely

$$e(\tilde{w}; \gamma) = \frac{1}{(1-\gamma)} \exp\left\{-\frac{1}{1-\gamma}(w_1 + w_2)\right\} I_0\left(\frac{2}{1-\gamma}\sqrt{\gamma w_1 w_2}\right) \quad (6)$$

where $0 \leq \gamma < 1$ and $I_0(a) = \sum_{k=1}^{\infty} \frac{a^{2k}}{4^k (k!)^2}$ is the modified Bessel function of the first kind of order zero. The value γ represents the Pearson's product-moment correlation and $\gamma = 0$ implies independence between the w_1 and w_2 . Coupling (5) with the density (6) we obtain a flexible specification of a bivariate quantile regression model that covers the case of independence between the components when $R(\phi)$ is an identity matrix and $\gamma = 0$.

3.2 Model-based clustering

Suppose that we observe data over n statistical units and, for simplicity, the same number $q \times T$ of values y_{ijt} , $j = 1, \dots, q$, $t = 1, \dots, T$, for each statistical unit i , $i = 1, \dots, n$. We collect the observations for each unit into the vector $\mathbf{y}_i = (\tilde{y}'_{i1}, \dots, \tilde{y}'_{iT})'$, with $\tilde{y}_{it} = (y_{i1t}, y_{i2t}, \dots, y_{iqt})'$, $t = 1, \dots, T$. The vectors \mathbf{y}_i are supposed to be independent. The whole dataset will be denoted by $\mathcal{Y} = (\mathbf{y}_1, \dots, \mathbf{y}_n)$.

Our statistical problem is to cluster the n statistical units, i.e. the sites, into the clusters $\{1, \dots, K\}$, $K \ll n$. We follow a mixture-model approach (McLachlan and Peel, 2000) to clustering according to which the cluster membership of the i th unit is represented by a latent random variable $c_i \in \{1, \dots, K\}$, where $c_i = k$ indicates if the i -th units belongs to the cluster k .

The memberships c_i , $i = 1, \dots, n$ are supposed independent and identically distributed variables with $\Pr(c_i = k) = \alpha_k$, $0 < \alpha_k < 1$, for all $k = 1, \dots, K$ and $\sum_{k=1}^K \alpha_k = 1$.

Given the cluster membership, c_i , observations for the i -th unit are generated by the mixture model specified hierarchically, namely

the observation level

$$\tilde{y}_{it}|c_i, \tilde{w}_{it} \sim \mathcal{N}_q(X_t' \tilde{\beta}_{c_i} + D_{c_i} \Theta \tilde{w}_{it}, W_{it} D_{c_i} \Sigma_{c_i} D_{c_i} W_{it}), \quad i = 1, \dots, n, \quad t = 1, \dots, T \quad (7)$$

conditionally independent distributed, with $\Sigma_{c_i} = \Omega R(\phi_{c_i}) \Omega$;

the latent process level

$$\tilde{w}_{it}|c_i \sim e(w, \gamma_{c_i}), \quad i = 1, \dots, n, \quad t = 1, \dots, T \quad (8)$$

conditionally independent distributed

$$c_i \sim \text{Multinomial}(1, \alpha), \quad i = 1, \dots, n \quad (9)$$

independent distributed. However, in presence of statistical units over a spatial domain, the incorporation of the spatial dependence in the model can be a feature that needs to be considered. Since the aim of our main application (see Section 5) was to discover sites at different environmental pressure, the incorporation of spatial dependence in the MCMC chain can eventually mask the classification by a common smoothing effect. Depending on the type of spatial domain and the analysed data, this dependence can be incorporated, for example, by considering a Markov random field as, for example, in (Gaetan *et al.*, 2017; Jiang and Serban, 2012) that takes into account the membership of the nearest neighbors.

3.3 Bayesian inference

We adopt a Bayesian approach to make inference on the model parameters. The inference is facilitated by the fact that the mixture model is specified hierarchically. Moreover the conjugacy of some priors leads to updates with simple and established methods for drawing from the full conditional distribution.

Prior distributions

A common choice for the prior distribution of α is $\alpha \sim \mathcal{D}(a_1, \dots, a_K)$, where \mathcal{D} is the Dirichlet distribution with parameters $a_1, \dots, a_K > 0$;

The choice $\tilde{\beta}_k \sim \mathcal{N}(b_k, P_k)$, $k = 1, \dots, K$, identically and independent distributed is still common and simplifies the simulation. Setting $b_k = 0$ and $P_k = aI$, for $a \gg 0$, leads to an improper prior.

Specification of the prior distributions for $\tilde{\sigma}_1, \dots, \tilde{\sigma}_k$ and $\tilde{\phi}_1, \dots, \tilde{\phi}_k$ is complicated by the complex requirement that the matrices $D_k \Omega R(\phi_k) \Omega D_k$ are non-negative definite matrices. We follow Barnard *et al.* (2000) and work by specifying prior for variances and correlation matrices. For the variances we put the priors $\sigma_{kj}^2 \sim \text{IG}(s_k, d_k)$, $k = 1, \dots, K$, $j = 1, \dots, q$, identically and independent distributed, where IG is the inverse gamma distribution, with the shape s_k and d_k scale parameters. Barnard *et al.* (2000) discussed the relative merits of choosing a prior for $\tilde{\phi}_1, \dots, \tilde{\phi}_k$ independent from $\tilde{\sigma}_1, \dots, \tilde{\sigma}_k$. In our example the choice is greatly simplified since we consider two variables, i.e. $q = 2$. In that case we can assign a uniform prior, $\phi_k \sim \mathcal{U}(-1, 1)$, for the correlation coefficient ϕ_k that are supposed identically and independent distributed.

For sampling from the posterior distribution we use a hybrid MCMC algorithm known as Metropolis-within-Gibbs algorithm (Robert and Casella, 2004, Chapter 10). As we show in the Supplementary material the conjugacy of some priors leads to updates with simple, established methods for drawing from the full conditional distribution. In other cases, we resort to Metropolis-Hastings to draw from some of the full conditional distributions.

4 A simulation study

In this simulation study we want to exemplify how the mixture model described in subsection (3.3) is able to cluster different bivariate temporal patterns encountered.

We will consider two experiments. A first experiment (Sim A) in which data are generated from a symmetric distributions and a second example (Sim B) in which data

come from asymmetric distributions. In both experiments we assume that $n = 300$ statistical units are split into three clusters of size $n_k = 100$, $k = 1, 2, 3$. For each unit we simulate a bivariate vector of length $T = 100$ $(y_{1t}, y_{2t})'$, $t = 1, \dots, T$ with time-varying marginal distributions.

Let z_{jt} , $j = 1, 2$, $t = 1, \dots, T$ a standardized Gaussian random variable and $m_{jk}(t)$ a positive function, $k = 1, 2, 3$, $j = 1, 2$, $t = 1, \dots, T$. We consider two setting for simulating $(y_{1t}, y_{2t})'$, namely

Sim A: $y_{jt} = m_{jk}(t) + \sqrt{m_{jk}(t)/5} z_{jt}$;

Sim B: $y_{jt} = G^{-1}(\Phi(z_{jt}); m_{jk}(t)/5, 5)$, where $\Phi(z)$ is the cumulative distribution function (CDF) of a standardized Gaussian random variable and $G^{-1}(u; a, b)$ is the inverse of the CDF of a Gamma random variable with mean ab and variance ab^2 .

The k value in $m_{jk}(t)$ varies between 1 and 3 depending on the cluster membership. Note that the marginal distributions were chosen in the way that the means and the marginal variances are equal in the two settings.

The temporal patterns in each cluster are led by the function $g(t; a, b) = a[2 + t/T + \exp\{-(t/T - b)^2/0.05\}]$, that is

Cluster	Component	
	$j = 1$	$j = 2$
$k = 1$	$m_{11}(t) = g(t; 1, 0.2)$	$m_{21}(t) = g(t; 1.5, 0.8)$
$k = 2$	$m_{12}(t) = g(t; 1, 0.5)$	$m_{22}(t) = g(t; 1.5, 0.2)$
$k = 3$	$m_{13}(t) = g(t; 1, 0.8)$	$m_{23}(t) = g(t; 1.5, 0.5)$

In order to asses the robustness of the procedure in the presence of serial dependence, we simulate the bivariate time series $(z_{1t}, z_{2t})'$ as $z_{jt} = v_{jt} + \theta v_{jt-1}$, $j = 1, 2$, where $(v_{1t}, v_{2t})' \sim \mathcal{N}(0, \Sigma_v)$, is a bivariate white noise, with $\Sigma_v = \frac{1}{1 + \theta^2} \begin{bmatrix} 1 & \rho \\ \rho & 1 \end{bmatrix}$, $-1 < \rho < 1$ and $-1 \leq \theta \leq 1$. By choosing different values for ρ and θ , different degrees of mutual and serial dependence are obtained.

In order to capture the temporal component of the bivariate vector in each cluster, we follow a regression spline approach. We choose a cubic B-spline basis, $b_1(t), \dots, b_m(t)$, with equally spaced knots over the range of time. For simplicity, we assume the same number m of basis functions for both component of \tilde{y}_t . The resulting matrix X_t of regressors in (7) is $X_t = \begin{pmatrix} x_t & 0 \\ 0 & x_t \end{pmatrix}'$ with $x_t = (b_1(t), \dots, b_m(t))'$.

We suggest this simple strategy for a preliminary selection of the number of basis m , namely

1. fix the value p_j , $j = 1, 2$ in (1) and get $\hat{\beta}_j$, $j = 1, 2$, for each time series $\{y_{jt}, t = 1, \dots, T\}$;
2. evaluate the AIC-like criterion $AIC_j(m) = \sum_{t=1}^T \rho_p(y_{jt} - x_t' \hat{\beta}_j) + 2m$;
3. repeat step 1 and 2 for each statistical unit i and obtain the value $AIC_j^{(i)}(m)$, $i = 1, \dots, n$;
4. find the value m that minimizes the overall value $\overline{AIC}(m) = \sum_{i=1}^n \sum_{j=1}^2 AIC_j^{(i)}(m)$.

In our simulation study we consider the median value for both time series, i.e. $p_1 = p_2 = 0.5$.

For estimating the model parameters ψ , we run the MCMC for 100 iterations as burn-in and 300 iterations for getting the posterior estimates. Inspection of trace plots suggests convergence of the parameters. From the clustering partitions sampled in the MCMC algorithm we obtain an estimate of the clustering structure by considering the posterior mode.

Both experiments are repeated 100 times. We compared the performance of our clustering method with three state-of art competitors:

- 1) Gaussian finite Mixture model Clustering (GMC): we cluster the data by means of a mixture of Gaussian linear regression models. For estimating the parameter we exploit the R package `flexmix` (Grün and Leisch, 2008) which is based on Expectation-Maximization algorithm;

- 2) Raw Data Clustering (RDC): each bivariate time series is stacked in one vector. Then the vectors are clustered by means of a Partitioning Around Medoids (PAM) algorithm, extracting three clusters;
- 3) CHaracteristic-based Clustering (CHC). A global measure describing the time series is obtained by applying summary indices about trend, seasonality, periodicity, serial correlation, skewness, kurtosis, chaos, nonlinearity and self-similarity (Wang *et al.*, 2006). The normalized indices or features extracted by using the R package `tsfeatures` (Hyndman *et al.*, 2020) are the inputs of the PAM algorithm.

We assess the power of our clustering algorithm in reconstructing the three clusters by comparing the level of agreement between the estimated partition and the true clustering using the Adjusted Rand Index, ARI, (Hubert and Arabie, 1985).

Table 1 presents the ARI values for the combinations of the correlation structures and the marginal distribution. From these results we note that our bivariate clustering algorithm performs very well at the combination, reporting low clustering performance for the more complex structure in particular for $\phi=0.5$; the simulation with the Gamma distribution (Sim B) appears more challenging in the clustering for all the considered methods. It is worth to note that our method reported a better classification than the Gaussian mixed linear regression model even in the case of normal marginal distribution, especially in presence of a serial dependence.

	ρ	θ	Our method	GMC	RDC	CHC
Sim A	0	0	1.00 (0.01)	0.94 (0.10)	1.00 (0.01)	0.82 (0.08)
	0.5	0	1.00 (0.01)	0.95 (0.09)	1.00 (0.01)	0.80 (0.05)
	0.5	1.0	0.91 (0.15)	0.57 (0.06)	0.75 (0.09)	0.57 (0.07)
Sim B	0	0	0.84 (0.12)	0.44 (0.05)	0.35 (0.11)	0.55 (0.09)
	0.5	0	0.56 (0.29)	0.44 (0.05)	0.35 (0.10)	0.58 (0.06)
	0.5	1.0	0.32 (0.22)	0.32 (0.05)	0.22 (0.10)	0.42 (0.06)

Table 1: Average (standard error between parentheses) Adjusted Rand Index values of 100 replications for each ρ and θ combination, marginal distribution and clustering method.

5 Clustering sites in the Gulf of Gabes

In this section we present the clustering results for the bivariate variables Chl-a concentration and the KD-490 levels, previously presented in Section 2. As reported in Figure 3 and 1, we can observe a strong and time-varying seasonal pattern, with the presence of a peak at the beginning of each year. For this reason we choose to model such monthly seasonality using sine and cosine functions (Eilers *et al.*, 2008). More precisely we suppose that for each variable y the quantile regression function is a function of the time t , such that

$$Q_p(y|t) = g_1(t) + g_2(t) \cos(\pi/6t) + g_3(t) \sin(\pi/6t)$$

where a possible overall trend is represented by a smooth function $g_1(t)$, while $g_2(t)$ and $g_3(t)$ are smooth functions that modulate the local amplitudes of the cosine and sine waves.

In order to have a good grade of flexibility, the three functions $g_1(t)$, $g_2(t)$, and $g_3(t)$ take the form of a regression on a cubic B-spline basis, $g_j(t) = \sum_{l=1}^{m_j} \beta_{jl} b_l(t)$ with equally spaced knots over the time interval $[1, 108]$. It is easy to see that the resulting model for the quantile function

$$Q_p(y|t) = \sum_{l=1}^{m_1} \beta_{1l} b_l(t) + \sum_{l=1}^{m_2} \beta_{2l} [b_l(t) \cos(\pi/6t)] + \sum_{l=1}^{m_3} \beta_{3l} [b_l(t) \sin(\pi/6t)]$$

can be written as a linear combination of covariates that fits with (5).

We carry out a preliminary data analysis in order to get the degree of smoothing and the number of clusters. This data analysis has been performed on the time series of Chl-a concentration since this variable displayed more heterogeneity in space and time.

To find the grade of smoothing and obtain a value for m_1 , m_2 , and m_3 , we follow the strategy outlined in Section 4 and we minimize the overall value $AIC(m)$, with $m = m_1 + m_2 + m_3$ by means of a median (i.e. $p = 0.5$) regression and considering a number of basis for each component from 3 to 6. The solution with $m_1 = 4$ for the trend (with the inclusion of a internal intercept) and $m_2 = m_3 = 3$ for the cyclical components

minimized the overall $AIC(m)$ and it was chosen in the following models.

For a fixed number of clusters K , we fit a bivariate model assuming that the two time series have a constant correlation coefficient across the clusters, i.e. $\phi = \phi_k$. Moreover a word of caution is in order of the estimate of the parameter γ in (6). Our experience with this dataset indicates that the data provide very little information on the parameter. For this reason we set the value $\gamma = 0.5$, giving the parameter ϕ the task of modulating the dependence between the two time series.

The results were obtained after 2500 Monte Carlo iterations using a burn-in of 300. The final membership c_i and the regression coefficients β_{jl} were estimated by means of the mode and mean of a-posteriori distribution, respectively.

We identify K minimizing an adapted version of the Deviance Information Criterion (DIC) (Spiegelhalter *et al.*, 2002) following Celeux *et al.* (2006). More precisely starting from the formulation named DIC_2 in that paper

$$DIC_2(\mathbf{y}, \psi) = -4\mathbb{E}_{\psi, c}[\log f(\mathbf{y} | \psi) | \mathbf{y}] + 2\log f(\mathbf{y}, \tilde{\psi}(\mathbf{y}))$$

where $\psi = (\beta, \sigma)'$ and $\tilde{\psi}(\mathbf{y})$ is the maximum a posteriori (MAP) estimates of ψ . Celeux *et al.* (2006) approximates DIC_2 by using the MCMC runs

$$\begin{aligned} CDIC(\mathbf{y}) = & -\frac{4}{m} \sum_{l=1}^m \sum_{i=1}^n \log \left\{ \sum_{k=1}^K \alpha_k^{(l)} f(\mathbf{y}_i | \psi_k^{(l)}) \right\} \\ & + 2 \sum_{i=1}^n \log \left\{ \frac{1}{m} \sum_{l=1}^m \sum_{k=1}^K \alpha_k^{(l)} f(\mathbf{y}_i | \psi_k^{(l)}) \right\} \end{aligned} \quad (10)$$

where $\psi_k^{(m)}$ and $\alpha_k^{(l)}$ are the results of the l -th MCMC iteration.

However, Formula 10 entails the evaluation of the bivariate density function $f(y_{i1t}, y_{i2t}; \psi_k)$ that arises from equation (5) by integrating out the random variable \tilde{w} . Since the bivariate density function cannot be derived in closed form, we proposed a composite version of the DIC index, called Composite-DIC (CDIC), pretending in (10) that

$f(y_{i1t}, y_{i2t}; \psi_k) = f(y_{i1t}; \psi_k) \cdot f(y_{i2t}; \psi_k)$ in the same spirit of Varin and Vidoni (2005).

We estimate the model for a range of different quantile combinations namely for different values of (p_1, p_2) . In particular we consider the combination of quantiles given by the pairs $(0.5, 0.5)$, $(0.9, 0.5)$, and $(0.9, 0.9)$ for Chl-a concentration and KD-490, respectively. While the quantile 0.5 can represent a robust estimate of the central tendency of the behaviour of each indicator, the quantile 0.9 is particularly important in ecology to evaluate the temporal trend towards the upper end of the distribution.

The values of CDIC for a number of clusters K which varies from 2 to 7 are reported in Table 2.

K	(0.5,0.5)	(0.9,0.5)	(0.9,0.9)
7	33.1	46.1	48.4
6	31.2	44.1	46.5
5	33.1	44.0	48.2
4	35.1	39.8	46.9
3	32.3	44.1	46.4
2	39.5	41.0	46.6

Table 2: CDIC ($\times 10^5$) values for K at different quantile pairs. The lowest CDIC value is highlighted in bold.

Considering the values reported in Table 2 we have chosen a number of clusters equal to $K = 6, 4$ and 3 for the quantile combination $(0.5, 0.5)$, $(0.9, 0.5)$ and $(0.9, 0.9)$, respectively. This result suggests a different number of clusters for each combination indicating a decreasing number to an increasing combination of quantile levels. Especially for the latest model, the proposed criterion suggests a classification based on three clusters.

In Figure 4 we present the spatial distribution of the clustering results for each combination of quantiles.

Across all the fitted quantile combinations, we observe that those regions around the islands (Kerkennah, Kneiss and Jerba) are clustered as having the highest average values of both Chl-a and KD-490 concentration which decreases as one moves towards the deep sea. These results are consistent with the results obtained from the univariate case as well as the findings of Katlane *et al.* (2012) who notes that, from multi-temporal turbidity

maps produced from Moderate Resolution Imaging Spectrometer (MODIS) for 2009, areas around the islands (Kerkenah, Kneiss and Jerba) and at the industrial port of Gannouch, were characterized by high turbidity variation, concentration of total suspended matter and Chl-a concentration.

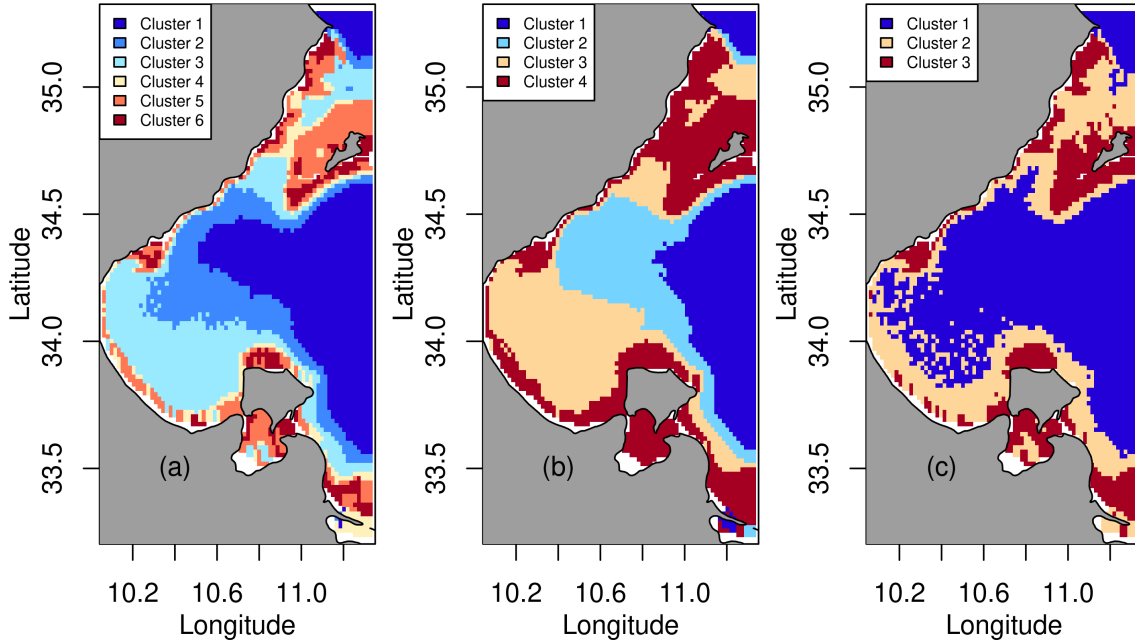


Figure 4: Spatial clustering results for Chl-a concentration and KD-490 index for each quantile combination (p_1, p_2) : (a) (0.5, 0.5) (left), (b) (0.9, 0.5) (middle), and (c) (0.9, 0.9) (right).

However, given the different meanings of Chl-a concentration and KD-490 levels at each quantile level, the spatial distribution between the three classifications appears different in particular looking at the northeastern area and the coastal zone. The classification with the quantile combination (0.5, 0.5) identifies a series of clusters that grades not only the most polluted area (Cluster 6) but also the coastal area (Cluster 4 and 5). The zone inside the Gulf is classified with the Cluster 3, while the remaining zone far from the coastal area and islands is covered by the Cluster 1 and 2. Otherwise, the classification with (0.9, 0.5) reported a similar spatial extension concerning the Cluster 1 and 2 which include zones starting from the offshore waters to the coastal area, while the coastal zone and the Sfax industrial area in the north-east are entirely covered by the Cluster 4; the Cluster 3 defines a transition zone. As reported in Figure 5, both the classification

performed with (0.5, 0.5) and (0.9, 0.5) report an increasing trend in the observed values of both Chl-a concentrations and KD-490 levels with the increase of the cluster label. The clustering performed with a quantile combination of 0.9 for Chl-a and 0.5 for KD-490 appears different with respect to the combination (0.5, 0.5) both in terms of marginal and spatial distribution: in the quantile combination (0.9, 0.5) Cluster 4 embraces an increased percentage of the seawater (24.6%) and it shows lower values with respect to the previous most impacted area (Cluster 6) obtained by the classification with quantiles (0.5, 0.5); low KD-490 and CHL-a concentrations are restricted for both the classifications ((0.5, 0.5) and (0.9, 0.5)) to the Cluster 1, 2 and Cluster 3. In the classification with (0.9, 0.9), the number of sites that belongs to Cluster 1 increases to 56.3%, and the relative spatial extension covers all the offshore water area. The Chl-a and KD-490 values in the Cluster 1 are low and widely separated from the Cluster 2 in comparison with the previous classifications. The sites belonging to Cluster 2 cover in a unique group all the sites previously classified as a transition zone or near the coastal area. The Cluster 3 identifies four marine areas heavily impacted by high measurements of Chl-a and KD-490: the seawater around Jerba, Sfax, Zarzis, and Sharqi Island.

In Figure 6 the temporal pattern for each recovered cluster is reported by plotting the estimated temporal component ($\hat{\beta}X$) within each identified group of Chl-a concentration and KD-490 level for each considered combination of quantiles. For all clusters a seasonal pattern is evident. Cluster labels and colors are the same as those in Figure 4 and they are ordered by increasing average Chl-a concentration. The classification with the quantile combination (0.5, 0.5) is guided by different average levels for Chl-a concentration and KD-490 level; however, especially for Chl-a concentration, the first three clusters (1, 2, and 3) look very close; the presence of a different classification for those groups is explained by a clear separation considering respective KD-490 levels: posing our attention to KD-490, the curves report different average values, but a similar seasonal component. In addition, the Cluster 6 reports an increasing trend at the beginning of the temporal window for both the indicators and a strong cyclical pattern in Chl-a concentration, while the seasonality

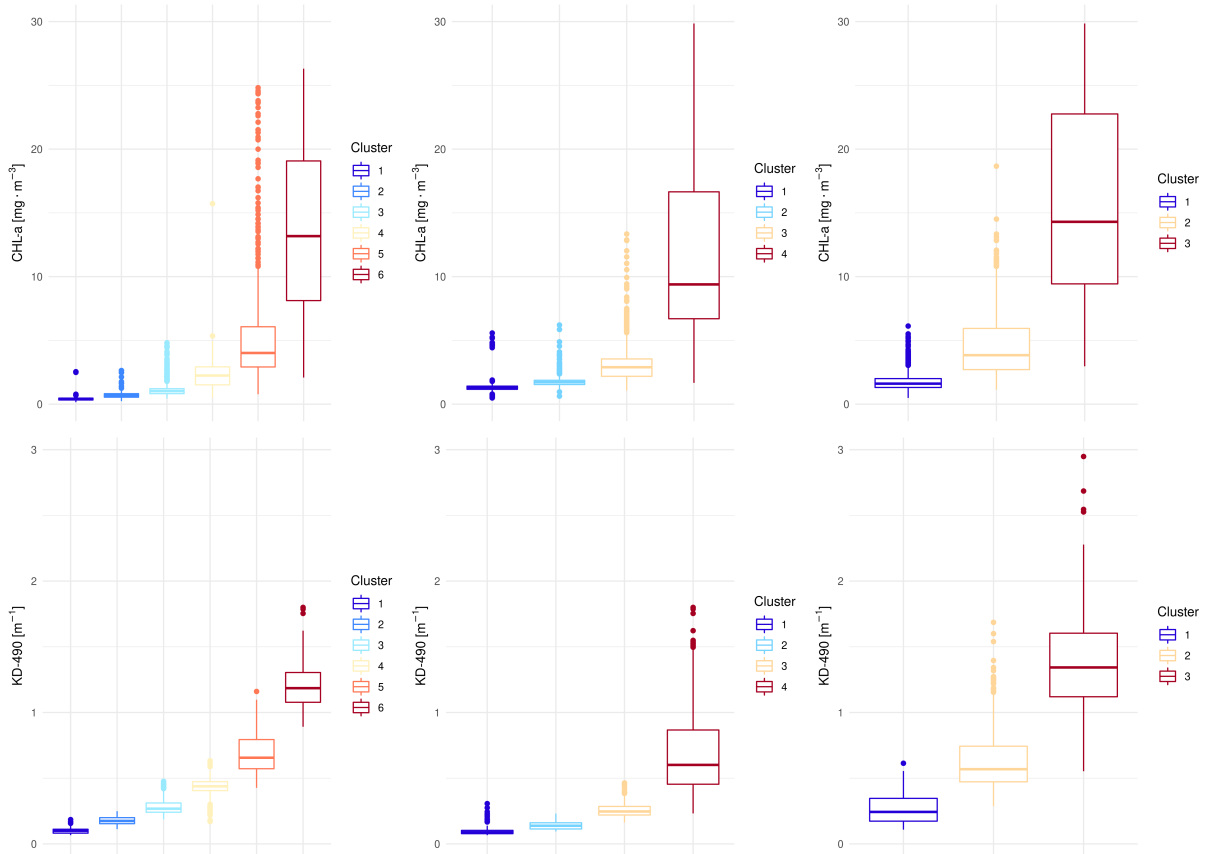


Figure 5: Distribution of reference quantile of Chl-a concentration and KD-490 index by Cluster for each quantile combination (p_1, p_2) : $(0.5, 0.5)$ (left), $(0.9, 0.5)$ (middle) and $(0.9, 0.9)$ (right).

is less evident for KD-490.

In Figure 5, considering the two boxplots in the middle related to the clustering results obtained with the quantile combination $(0.9 - 0.5)$, the analysis of the differences between the estimated groups reports as sites belonging to Cluster 4 are those with the highest values and a well separated from the other clusters for both the indicators. Otherwise, Cluster 1, 2, and 3 are close together and the difference is mainly due to a different average trend, a greater intra-season amplitude for Cluster 3, and a late seasonal peak for Cluster 1, as denoted by the Figure 6. In addition, the temporal trend appears to be different from the previous classification: only for KD-490, Cluster 4 reports an evident increasing trend, followed by a stabilization and a slight decrease after the year 2008. Taking into account the clustering results obtained by the quantiles $(0.9 - 0.9)$, all clusters are well separated.

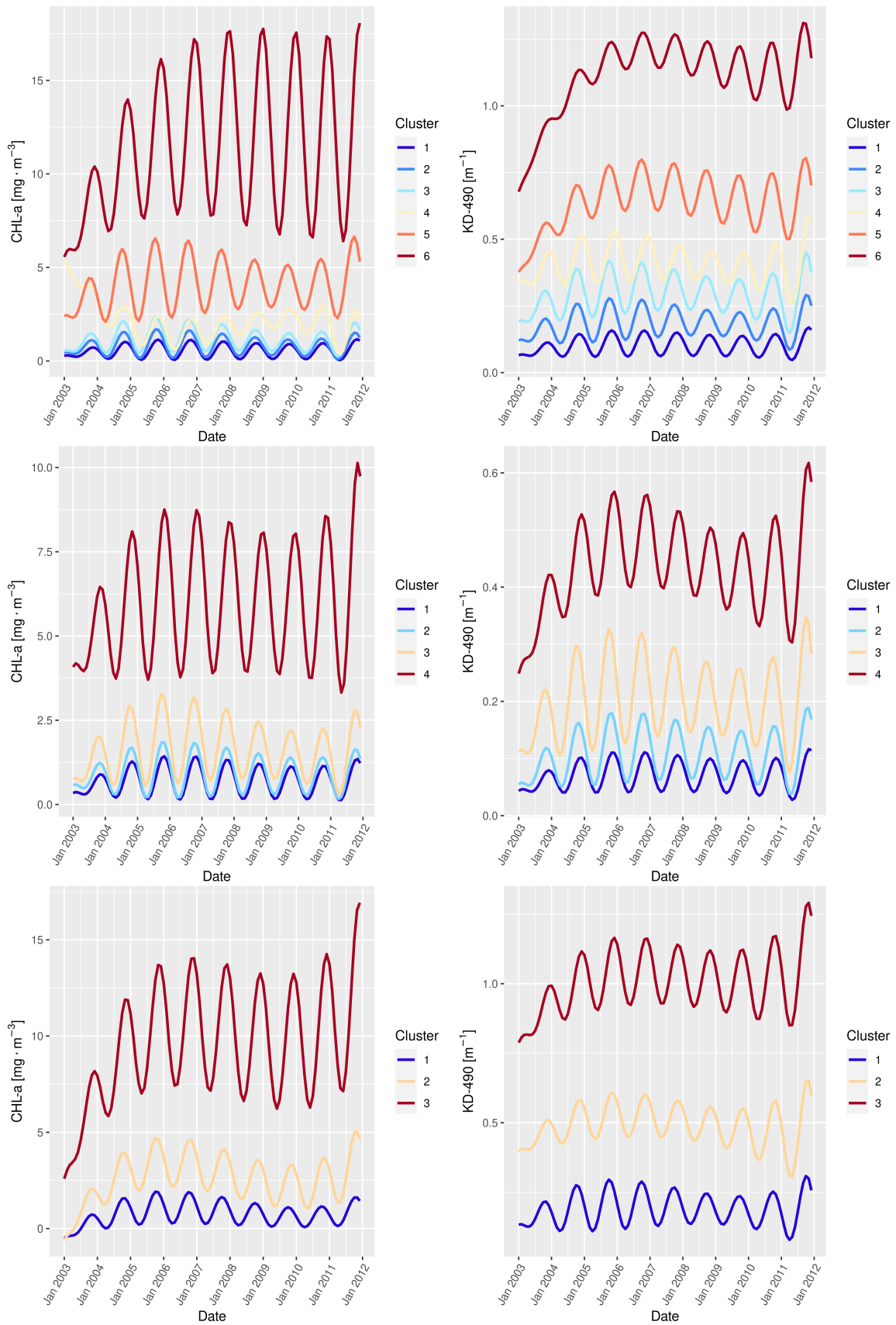


Figure 6: Estimated temporal component for Chl-a concentration and KD-490 index at quantile combinations (p_1, p_2) : (0.5, 0.5) (top), (0.9, 0.5) (middle), and (0.9, 0.9) (bottom).

Cluster 1 reports the lowest values of limited signal amplitude, while the Cluster 2 shows a stable and cyclical behaviour, far from other groups. Cluster 3 exhibits an increasing trend and a stabilization after 2006 only, for Chl-a concentrations, while the KD-490 trend is stable and cyclical.

6 Discussion

In this paper we have proposed a new model-based clustering technique, that is capable of handling asymmetric clusters with the presence of outliers as well as considering different quantile levels of the observed data. Our clustering strategy is based on the finite mixture model theory where each component of the mixture of AL distributions, which constitutes the density of bivariate random variables which are potentially correlated, is assumed to represent a cluster with the skewness parameter of AL distribution being used to directly model the quantiles of interest. Therefore our proposed technique adds to the richness of the recent burgeoning of non-Gaussian approaches to model-based clustering.

As the AL distribution provides a direct link between the maximum likelihood theory and minimization of a quantile regression check loss function (see Koenker and Machado, 1999; Yu and Moyeed, 2001), we estimate the cluster-specific parameters, which are the parameters of the mixing AL distribution, through a Bayesian approach.

In our simulation experiment we considered three clusters two of which are not distinctively different from each other. However, an evaluation of the power of our proposed algorithm in reconstructing the three groups indicates a good performance with respect to the other competitive methods. We have applied the procedure to time series observed from GlobColour data related to Chlorophyll type-a concentrations and KD-490 levels in order to identify homogeneous areas in the Gulf of Gabes with respect to the temporal behavior of these water indicators by means of a seasonal modulation model. We defined clusters that are similar by different combination of quantiles of the two indicators. It is important to note that as different choice of quantiles implies changes in the clustering,

this method may be particularly suitable for defining areas at different risk when considering two indicators at different quantile levels. More particularly, important features are of absolute interest in environmental sciences and ecology (Schmidt *et al.*, 2012). In addition, the use of a model matrix based on a seasonal modulation model is evoked by the periodic behaviour of time series in our environmental application; however other flexible specifications can also be adopted.

We note that both Chl-a concentrations and KD-490 levels are affected by several spatially varying factors. Potential consequences of non incorporating in the model information related to the spatial domain may result from misclassification to lower predictive ability; the spatial dependence can be particularly helpful in presence of high percentage of missing data. In this case potential future development of the modelling approach can be the incorporation of spatial dependence among the probabilities of membership as in Jiang and Serban (2012) or Gaetan *et al.* (2017). Another possible extension is to perform clustering at multiple quantiles instead of fixing the levels of quantiles. However, caution has to be taken in this case to avoid the issue of crossing quantiles.

References

- Alikas, K., Kangro, K., Randoja, R., Philipson, P., Asuküll, E., Pisek, J., and Reinart, A. (2015). Satellite-based products for monitoring optically complex inland waters in support of EU Water Framework Directive. *International Journal of Remote Sensing*, **36**, 4446–4468.
- Aloulou, F., ElEuch, B., and Kallel, M. (2012). Benthic foraminiferal assemblages as pollution proxies in the northern coast of Gabes Gulf, Tunisia. *Environmental Monitoring and Assessment*, **184**, 777–795.
- Alvera-Azcárate, A., Sirjacobs, D., Barth, A., and Beckers, J.-M. (2012). Outlier detection in satellite data using spatial coherence. *Remote Sensing of Environment*, **119**, 84–91.

- Ayadi, N., Aloulou, F., and Bouzid, J. (2015). Assessment of contaminated sediment by phosphate fertilizer industrial waste using pollution indices and statistical techniques in the Gulf of Gabes (Tunisia). *Arabian Journal of Geosciences*, **8**, 1755–1767.
- Barbosa, S., Scotto, M., and Alonso, A. (2011). Summarising changes in air temperature over Central Europe by quantile regression and clustering. *Natural Hazards and Earth System Sciences*, **11**, 3227–3233.
- Barnard, J., McCulloch, R., and Meng, X.-L. (2000). Modeling covariance matrices in terms of standard deviations and correlations, with application to shrinkage. *Statistica Sinica*, **10**, 1281–1311.
- Benoit, D. F. and Van den Poel, D. (2012). Binary quantile regression: a Bayesian approach based on the asymmetric Laplace distribution. *Journal of Applied Econometrics*, **27**, 1174–1188.
- Benoit, D. F., Alhamzawi, R., and Yu, K. (2013). Bayesian lasso binary quantile regression. *Computational Statistics*, **28**, 2861–2873.
- Cazelles, B., Chavez, M., Berteaux, D., Ménard, F., Vik, J. O., Jenouvrier, S., and Stenseth, N. C. (2008). Wavelet analysis of ecological time series. *Oecologia*, **156**, 287–304.
- Celeux, G., Forbes, F., Robert, C. P., and Titterton, D. M. (2006). Deviance information criteria for missing data models. *Bayesian Analysis*, **1**, 651–673.
- Dabuleviciene, T., Vaiciute, D., and Kozlov, I. E. (2020). Chlorophyll-a variability during upwelling events in the South-Eastern Baltic Sea and in the Curonian Lagoon from satellite observations. *Remote Sensing*, **12**, 3661.
- Directive, E. W. F. *et al.* (2000). The European parliament and of the council. *Water Framework Directive (2000/60/EC)*, *OJL*, **327**, 1–73.
- Downton, F. (1970). Bivariate exponential distributions in reliability theory. *Journal of the Royal Statistical Society: Series B*, **32**, 408–417.

- Eilers, P. H., Gampe, J., Marx, B. D., and Rau, R. (2008). Modulation models for seasonal time series and incidence tables. *Statistics in Medicine*, **27**, 3430–3441.
- El Kateb, A., Stalder, C., Neururer, C., Pisapia, C., and Spezzaferri, S. (2016). Correlation between pollution and decline of Scleractinian Cladocora Caespitosa (Linnaeus, 1758) in the Gulf of Gabes. *Heliyon*, **2**, e00195.
- El Zrelli, R., Courjault-Radé, P., Rabaoui, L., Daghbouj, N., Mansour, L., Balti, R., Castet, S., Attia, F., Michel, S., and Bejaoui, N. (2017). Biomonitoring of coastal pollution in the Gulf of Gabes (se, Tunisia): use of *Posidonia oceanica* seagrass as a bioindicator and its mat as an archive of coastal metallic contamination. *Environmental Science and Pollution Research*, **24**, 22214–22225.
- El Zrelli, R., Rabaoui, L., Alaya, M. B., Daghbouj, N., Castet, S., Besson, P., Michel, S., Bejaoui, N., and Courjault-Radé, P. (2018). Seawater quality assessment and identification of pollution sources along the central coastal area of Gabes Gulf (se Tunisia): evidence of industrial impact and implications for marine environment protection. *Marine Pollution Bulletin*, **127**, 445–452.
- Finazzi, F., Haggarty, R., Miller, C., Scott, M., and Fasso, A. (2015). A comparison of clustering approaches for the study of the temporal coherence of multiple time series. *Stochastic Environmental Research and Risk Assessment*, **29**, 463–475.
- Fourati, R., Tedetti, M., Guigue, C., Goutx, M., Zaghden, H., Sayadi, S., and Elleuch, B. (2018). Natural and anthropogenic particulate-bound aliphatic and polycyclic aromatic hydrocarbons in surface waters of the Gulf of Gabès (Tunisia, southern Mediterranean sea). *Environmental Science and Pollution Research*, **25**, 2476–2494.
- Gaetan, C., Girardi, P., Pastres, R., Mangin, A., *et al.* (2016). Clustering chlorophyll-a satellite data using quantiles. *Annals of Applied Statistics*, **10**, 964–988.
- Gaetan, C., Girardi, P., and Pastres, R. (2017). Spatial clustering of curves with an application of satellite data. *Spatial Statistics*, **20**, 110–124.

- Giraldo, R., Delicado, P., and Mateu, J. (2012). Hierarchical clustering of spatially correlated functional data. *Statistica Neerlandica*, **66**, 403–421.
- Grün, B. and Leisch, F. (2008). FlexMix version 2: Finite mixtures with concomitant variables and varying and constant parameters. *Journal of Statistical Software*, **28**, 1–35.
- Haggarty, R., Miller, C., and Scott, E. (2015). Spatially weighted functional clustering of river network data. *Journal of the Royal Statistical Society. Series C, Applied Statistics*, **64**, 491–506.
- Hamza, A. and El Abed, A. (1994). Les eaux colorées dans le golfe de Gabès: bilan de six ans de surveillance (1989-1994). *Bulletin de l’Institut National des Sciences et Technologies de la Mer*, **21**, 66–72.
- Hubert, L. and Arabie, P. (1985). Comparing partitions. *Journal of Classification*, **2**, 193–218.
- Hyndman, R., Kang, Y., Montero-Manso, P., Talagala, T., Wang, E., Yang, Y., and O’Hara-Wild, M. (2020). *tsfeatures: Time Series Feature Extraction*. R package version 1.0.2.
- Jiang, H. and Serban, N. (2012). Clustering random curves under spatial interdependence with application to service accessibility. *Technometrics*, **54**, 108–119.
- Jorgensen, B. (1982). *Statistical Properties of the Generalized Inverse Gaussian Distribution*. Springer-Verlag, New York.
- Katlane, R., DUPOUY, C., and Zargouni, F. (2012). Chlorophyll and turbidity concentrations deduced from MODIS as an index of water quality of the Gulf of Gabes in 2009. In AUF, editor, *Téledétection 11, 1*, Téledétection, pages 265–273. CNRS & Campus Spatial Univ. Paris Diderot VII.
- Koenker, R. and Bassett, G. (1978). Regression quantiles. *Econometrica*, **46**, 33–50.

- Koenker, R. and Machado, J. A. (1999). Goodness of fit and related inference processes for quantile regression. *Journal of the American Statistical Association*, **94**, 1296–1310.
- Kotz, S., N., B., and Johnson, N. (2000). *Continuous Multivariate Distributions. Volume 1: Models and Applications*. Wiley, New York.
- Kotz, S., Kozubowski, T., and Podgorski, K. (2001). *The Laplace Distribution and Generalizations: A Revisit with Applications to Communications, Economics, Engineering, and Finance*. Springer, New York.
- Li, H., Deng, X., Dolloff, C., and Smith, E. (2016). Bivariate functional data clustering: grouping streams based on a varying coefficient model of the stream water and air temperature relationship. *Environmetrics*, **27**, 15–26.
- Liechty, J. C., Liechty, M. W., and Müller, P. (2004). Bayesian correlation estimation. *Biometrika*, **91**, 1–14.
- McLachlan, G. J. and Peel, D. (2000). *Finite Mixture Models*. Wiley, New York.
- Monteiro, A., Carvalho, A., Ribeiro, I., Scotto, M., Barbosa, S., Alonso, A., Baldasano, J., Pay, M., Miranda, A., and Borrego, C. (2012). Trends in ozone concentrations in the Iberian Peninsula by quantile regression and clustering. *Atmospheric Environment*, **56**, 184–193.
- Pardo, I., Gómez-Rodríguez, C., Wasson, J.-G., Owen, R., van de Bund, W., Kelly, M., Bennett, C., Birk, S., Buffagni, A., Erba, S., *et al.* (2012). The European reference condition concept: a scientific and technical approach to identify minimally-impacted river ecosystems. *Science of the Total Environment*, **420**, 33–42.
- Petrella, L. and Raponi, V. (2019). Joint estimation of conditional quantiles in multivariate linear regression models with an application to financial distress. *Journal of Multivariate Analysis*, **173**, 70–84.

- Poikāne, S., Alves, M. H., Argillier, C., Van den Berg, M., Buzzi, F., Hoehn, E., De Hoyos, C., Karottki, I., Laplace-Treyture, C., Solheim, A. L., *et al.* (2010). Defining chlorophyll-a reference conditions in European lakes. *Environmental Management*, **45**, 1286–1298.
- Rabaoui, L., Balti, R., Zrelli, R., and Tlig-Zouari, S. (2013). Assessment of heavy metals pollution in the Gulf of Gabes (Tunisia) using four mollusk species. *Mediterranean Marine Science*, **15**, 45–58.
- Robert, C. P. and Casella, G. (2004). *Monte Carlo Statistical Methods*. Springer, New York.
- Schmidt, T. S., Clements, W. H., and Cade, B. S. (2012). Estimating risks to aquatic life using quantile regression. *Freshwater Science*, **31**, 709–723.
- Shi, K., Li, Y., Li, L., Lu, H., Song, K., Liu, Z., Xu, Y., and Li, Z. (2013). Remote chlorophyll-a estimates for inland waters based on a cluster-based classification. *Science of the Total Environment*, **444**, 1–15.
- Sottile, G. and Adelfio, G. (2019). Clusters of effects curves in quantile regression models. *Computational Statistics*, **34**, 551–569.
- Spiegelhalter, D. J., Best, N. G., Carlin, B. P., and Van Der Linde, A. (2002). Bayesian measures of model complexity and fit. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, **64**, 583–639.
- Stafoggia, M., Schwartz, J., Badaloni, C., Bellander, T., Alessandrini, E., Cattani, G., De’Donato, F., Gaeta, A., Leone, G., Lyapustin, A., *et al.* (2017). Estimation of daily PM10 concentrations in Italy (2006–2012) using finely resolved satellite data, land use variables and meteorology. *Environment International*, **99**, 234–244.
- Varin, C. and Vidoni, P. (2005). A note on composite likelihood inference and model selection. *Biometrika*, **92**, 519–528.

- Waldmann, E. and Kneib, T. (2015). Bayesian bivariate quantile regression. *Statistical Modelling*, **15**, 326–344.
- Wang, X., Smith, K., and Hyndman, R. (2006). Characteristic-based clustering for time series data. *Data Mining and Knowledge Discovery*, **13**, 335–364.
- Yang, C., Ye, H., and Tang, S. (2020). Seasonal variability of diffuse attenuation coefficient in the Pearl river estuary from long-term remote sensing imagery. *Remote Sensing*, **12**, 2269.
- Yu, K. and Moyeed, R. A. (2001). Bayesian quantile regression. *Statistics & Probability Letters*, **54**, 437–447.
- Zaghden, H., Kallel, M., Elleuch, B., Oudot, J., Saliot, A., and Sayadi, S. (2014). Evaluation of hydrocarbon pollution in marine sediments of Sfax coastal areas from the Gabes Gulf of Tunisia, Mediterranean Sea. *Environmental Earth Sciences*, **72**, 1073–1082.
- Zhang, Y., Wang, H. J., and Zhu, Z. (2019). Quantile-regression-based clustering for panel data. *Journal of Econometrics*, **213**, 54–67.

Supplementary material for

“Clustering of bivariate satellite time series: a quantile approach”

We observe the set $\mathcal{Y} = \{\mathbf{y}_i, i = 1, \dots, n\}$ of n vectors of independent observations, with $\mathbf{y}_i = (\tilde{y}'_{i1}, \dots, \tilde{y}'_{iT})'$, and $\tilde{y}_{it} = (y_{i1t}, y_{i2t}, \dots, y_{iqt})'$, $t = 1, \dots, T$ from the hierarchical model

the observation level

$$\tilde{y}_{it}|c_i, \tilde{w}_{it} \sim \mathcal{N}_q(X'_t \tilde{\beta}_{c_i} + D_{c_i} \Theta \tilde{w}_{it}, W_{it} D_{c_i} \Sigma_{c_i} D_{c_i} W_{it}), \quad i = 1, \dots, n, \quad t = 1, \dots, T$$

conditionally independent distributed, with $\Sigma_{c_i} = \Omega R(\tilde{\phi}_{c_i}) \Omega$;

the latent process level

$$\tilde{w}_{it}|c_i \sim e(w, \gamma_{c_i}) \quad i = 1, \dots, n, \quad t = 1, \dots, T$$

conditionally independent distributed

$$c_i \sim \text{Multinomial}(1, \alpha), \quad i = 1, \dots, n$$

independent distributed.

The prior distribution for the the model parameters $\boldsymbol{\psi} = (\alpha', \boldsymbol{\beta}', \boldsymbol{\gamma}', \boldsymbol{\sigma}', \boldsymbol{\phi}')$ with $\alpha = (\alpha_1, \dots, \alpha_K)'$, $\boldsymbol{\beta} = (\tilde{\beta}'_1, \dots, \tilde{\beta}'_K)'$, $\boldsymbol{\gamma} = (\gamma'_1, \dots, \gamma'_K)'$, $\boldsymbol{\sigma} = (\tilde{\sigma}'_1, \dots, \tilde{\sigma}'_K)'$ and $\boldsymbol{\phi} = (\tilde{\phi}'_1, \dots, \tilde{\phi}'_K)'$. are given by

$$\pi(\boldsymbol{\beta}) = \prod_{k=1}^K \pi(\tilde{\beta}_k) \propto \prod_{k=1}^K |P_k|^{-1/2} \exp \left\{ -\frac{1}{2} (\tilde{\beta}_k - b_k)' P_k^{-1} (\tilde{\beta}_k - b_k) \right\};$$

$$\pi(\boldsymbol{\sigma}) = \prod_{k=1}^K \pi(\tilde{\sigma}_k) \propto \prod_{k=1}^K \prod_{j=1}^q \sigma_{kj}^{2-s_k-1} \exp(-d_k / \sigma_{kj}^2);$$

$$\pi(\boldsymbol{\phi}) = \prod_{k=1}^K \pi(\tilde{\phi}_k) \text{ and } \pi(\boldsymbol{\gamma}) = \prod_{k=1}^K \pi(\gamma_k). \text{ More details on the analytical form of } \pi(\tilde{\phi}_k)$$

and $\pi(\gamma_k)$ will be given later.

$$\pi(\alpha) = \frac{\Gamma(a_1 + \dots + a_K)}{\Gamma(a_1) \dots \Gamma(a_K)} \prod_{k=1}^K \alpha_k^{a_k - 1} \text{ with parameters } a_1, \dots, a_K > 0$$

Finally we denote the set $\mathcal{W} = (\mathbf{w}'_1, \dots, \mathbf{w}'_n)$, with latent vectors $\mathbf{w}_i = (\tilde{w}'_{i1}, \dots, \tilde{w}'_{iT})'$, $\tilde{w}_{it} = (w_{i1t}, y_{i2t}, \dots, w_{iqt})'$, $t = 1, \dots, T$ the (latent) vector of the cluster memberships $\mathcal{C} = (c_1, \dots, c_n)$

We now proceed to deriving the conditional distribution of the model parameters including also the conditional for the latent quantities \mathcal{W} and \mathcal{C} .

Applying the Bayes theorem we obtain the posterior distribution as

$$\begin{aligned} \pi(\alpha, \boldsymbol{\beta}, \boldsymbol{\gamma}, \boldsymbol{\sigma}, \boldsymbol{\phi}, \mathcal{C}, \mathcal{W} | \mathcal{Y}) &\propto \prod_{i=1}^n f(\mathbf{y}_i | \alpha, \tilde{\boldsymbol{\beta}}_{c_i}, \tilde{\boldsymbol{\sigma}}_{c_i}, \tilde{\boldsymbol{\phi}}_{c_i}, \mathbf{w}_i, c_i) \times \\ &\prod_{i=1}^n f(\mathbf{w}_i | \gamma_{c_i}) \times f(c_i | \alpha) \times \\ &\prod_{k=1}^K \pi(\tilde{\boldsymbol{\beta}}_k) \times \pi(\tilde{\boldsymbol{\sigma}}_k) \times \pi(\tilde{\boldsymbol{\phi}}_k) \times \pi(\boldsymbol{\gamma}_k) \times \pi(\alpha) \quad (\text{S.1}) \end{aligned}$$

where

$$\begin{aligned} f(\mathbf{y}_i | \alpha, \tilde{\boldsymbol{\beta}}_{c_i}, \tilde{\boldsymbol{\sigma}}_{c_i}, \tilde{\boldsymbol{\phi}}_{c_i}, \mathbf{w}_i, c_i) &\propto \prod_{t=1}^T |W_{it} D_{c_i} \Sigma_{c_i} D_{c_i} W_{it}|^{-1/2} \exp \left\{ -\frac{1}{2} (\tilde{\mathbf{y}}_{it} - X'_t \tilde{\boldsymbol{\beta}}_{c_i} - D_{c_i} \boldsymbol{\Theta} \tilde{\mathbf{w}}_{it})' \times \right. \\ &\left. [W_{it} D_{c_i} \Sigma_{c_i} D_{c_i} W_{it}]^{-1} (\tilde{\mathbf{y}}_{it} - X'_t \tilde{\boldsymbol{\beta}}_{c_i} - D_{c_i} \boldsymbol{\Theta} \tilde{\mathbf{w}}_{it}) \right\}; \end{aligned}$$

Full conditional for α

To derive the posterior density of the probability vector α , we first note that $\pi(\alpha, |\boldsymbol{\beta}, \boldsymbol{\gamma}, \boldsymbol{\sigma}, \boldsymbol{\phi}, \mathcal{C}, \mathcal{W}, \mathcal{Y}) \propto \pi(\alpha | \mathcal{C})$ and then

$$\pi(\alpha | \boldsymbol{\beta}, \boldsymbol{\sigma}, \boldsymbol{\gamma}, \boldsymbol{\phi}, \mathcal{C}, \mathcal{W}, \mathcal{Y}) \propto \prod_{i=1}^n f(c_i | \alpha) \pi(\alpha) \propto \prod_{k=1}^K \alpha_k^{a_k + \sum_{i=1}^n I(c_i = k) - 1}$$

Therefore $\pi(\alpha | \boldsymbol{\beta}, \boldsymbol{\gamma}, \boldsymbol{\sigma}, \boldsymbol{\phi}, \mathcal{C}, \mathcal{W}, \mathcal{Y})$ is the density of a Dirichlet distribution with parameters $a_k + \sum_{i=1}^n I(c_i = k) - 1$, $k = 1, \dots, K$.

Full conditional for $\boldsymbol{\beta}$

We have

$$\pi(\boldsymbol{\beta} | \alpha, \boldsymbol{\gamma}, \boldsymbol{\sigma}, \mathcal{C}, \mathcal{W}, \mathcal{Y}) \propto \prod_{k=1}^K \pi(\tilde{\boldsymbol{\beta}}_k | \tilde{\boldsymbol{\sigma}}_k, \tilde{\boldsymbol{\phi}}_k, \mathcal{C}, \mathcal{W}, \mathcal{Y})$$

where

$$\pi(\tilde{\beta}_k | \tilde{\sigma}_k, \tilde{\phi}_k, \mathcal{C}, \mathcal{W}, \mathcal{Y}) \propto \exp \left\{ -\frac{1}{2} (\tilde{\beta}_k - b_k)' P_k^{-1} (\tilde{\beta}_k - b_k) \right\} \times \prod_{i=1}^n \left[\exp \left\{ -\frac{1}{2} \sum_{t=1}^T (\tilde{u}_{it} - X_t' \tilde{\beta}_{c_i})' S_{c_i, t}^{-1} (\tilde{u}_{it} - X_t' \tilde{\beta}_{c_i}) \right\} \right]^{I(c_i=k)},$$

with $\tilde{u}_{it} = (u_{i1t}, \dots, u_{iqt})' = \tilde{y}_{it} - D_{c_i} \Theta \tilde{w}_{it}$ and $S_{c_i, t} = W_{it} D_{c_i} \Sigma_{c_i} D_{c_i} W_{it}$

The previous formula can be further elaborated, namely

$$\pi(\tilde{\beta}_k | \tilde{\sigma}_k, \tilde{\phi}_k, \mathcal{C}, \mathcal{W}, \mathcal{Y}) \propto \exp \left\{ -\frac{1}{2} \left[\tilde{\beta}_k' \left(P_k^{-1} + \sum_{i=1}^n I(c_i = k) \tilde{X}' S_{c_i}^{-1} \tilde{X} \right) \tilde{\beta}_k - 2 \tilde{\beta}_k' \left(P_k^{-1} b_k + \sum_{i=1}^n I(c_i = k) \tilde{X}' S_{c_i}^{-1} \mathbf{u}_i \right) \right] \right\}$$

where S_{c_i} is a block diagonal matrix with entries $S_{c_i, t}$, $t = 1, \dots, T$, $\tilde{X} = [X_1, \dots, X_T]'$ is the matrix of covariates and $\mathbf{u}_i = (\tilde{u}_{i1}', \dots, \tilde{u}_{iT}')'$.

Therefore $\pi(\tilde{\beta}_k | \tilde{\sigma}_k, \tilde{\phi}_k, \mathcal{C}, \mathcal{W}, \mathcal{Y})$ is the density, up to a normalizing constant, of a multivariate Gaussian vector with vector mean

$$\bar{m}_k = \bar{S}_k^{-1} \left(P_k^{-1} b_k + \sum_{i=1}^n I(c_i = k) X' S_{c_i}^{-1} \mathbf{u}_i \right)$$

and covariance matrix

$$\bar{S}_k = \left(P_k^{-1} + \sum_{i=1}^n I(c_i = k) X' S_{c_i}^{-1} X \right)^{-1}.$$

Full conditional for σ

We have

$$\pi(\sigma | \alpha, \beta, \gamma, \phi, \mathcal{C}, \mathcal{W}, \mathcal{Y}) \propto \prod_{k=1}^K \pi(\tilde{\sigma}_k | \tilde{\beta}_k, \tilde{\phi}_k, \mathcal{C}, \mathcal{W}, \mathcal{Y})$$

where

$$\begin{aligned} \pi(\tilde{\sigma}_k | \tilde{\beta}_k, \tilde{\phi}_k, \mathcal{C}, \mathcal{W}, \mathcal{Y}) &\propto \prod_{j=1}^q [\sigma_{kj}^2]^{-(s_j+1)} \exp\left(-\frac{d_j}{\sigma_{kj}^2}\right) \times \\ &\prod_{i=1}^n \left\{ |\tilde{S}_{c_i}|^{-1/2} \exp\left[-\frac{1}{2} \sum_{t=1}^T \left(\tilde{u}_{it} - X'_t \tilde{\beta}_{c_i}\right)' S_{c_i,t}^{-1} \left(\tilde{u}_{it} - X'_t \tilde{\beta}_{c_i}\right)\right] \right\}^{I(c_i=k)}, \end{aligned}$$

However, it is difficult to calculate the normalizing constant for this density and sample from the conditional distribution. As a result we introduce a Metropolis-Hastings (M-H) step in our simulation algorithm. The proposals in the M-H step are the following.

We pretend that $\tilde{\phi}_k = 0$ and in this case

$$\pi(\tilde{\sigma}_k | \tilde{\beta}_k, 0, \mathcal{C}, \mathcal{W}, \mathcal{Y}) = \prod_{j=1}^q \pi(\sigma_{kj} | \tilde{\beta}_k, 0, \mathcal{C}, \mathcal{W}, \mathcal{Y})$$

with

$$\begin{aligned} \pi(\sigma_{kj} | \tilde{\beta}_k, 0, \mathcal{C}, \mathcal{W}, \mathcal{Y}) &\propto [\sigma_{kj}^2]^{-(s_j+1)} \exp\left(-\frac{d_j}{\sigma_{kj}^2}\right) \times \\ &\prod_{i=1}^n \left\{ (\sigma_{c_i,j}^2)^{-T/2} \exp\left[-\frac{1}{2\sigma_{c_i,j}^2 \omega^2} \sum_{t=1}^T \frac{(u_{ijt} - X'_{jt} \beta_{c_i})^2}{w_{ijt}}\right] \times \right. \\ &\left. (\sigma_{c_i,j}^2)^{-T} \exp\left(-\frac{1}{\sigma_{c_i,j}^2} \sum_{t=1}^T w_{ijt}\right) \right\}^{I(c_i=k)} \\ &\propto (\sigma_{kj}^2)^{-(2s_j+3n_k T)/2-1} \times \\ &\exp\left\{-\frac{1}{\sigma_{kj}^2} \left[d_j + \sum_{i=1}^n \left(\sum_{t=1}^T w_{ijt} + \frac{1}{2\omega^2} \sum_{t=1}^T \frac{(u_{ijt} - X'_{jt} \beta_{c_i})^2}{w_{ijt}} \right) I(c_i = k) \right] \right\} \end{aligned}$$

Here X_{jt} is the j -th row of the matrix X_t and $n_k = \sum_{i=1}^n I(c_i = k)$, the number of vectors in \mathcal{Y} with membership k .

Thus $\pi(\sigma_{kj} | \tilde{\beta}_k, 0, \mathcal{C}, \mathcal{W}, \mathcal{Y})$ is the distribution of an inverse Gamma random variable,

InvGamma(a_{kj}, b_{kj}), with shape and scale parameter

$$\begin{aligned} a_{kj} &= (2s_j + 3n_k T)/2 \\ b_{kj} &= \left[d_j + \sum_{i=1}^n \left(\sum_{t=1}^T w_{ijt} + \frac{1}{2\omega^2} \sum_{t=1}^T \frac{(u_{ijt} - X'_{jt}\beta_{c_i})^2}{w_{ijt}} \right) I(c_i = k) \right] \end{aligned} \quad (\text{S.2})$$

Note that random samples from inverse Gamma distribution can be drawn from a Gamma distribution exploiting the fact that if $G \sim \text{Gamma}(a, 1/b)$ then $1/G \sim \text{InvGamma}(a, b)$.

The proposal in M-H step is accomplished by sampling independent values $\sigma_{kj} \sim \text{InvGamma}(a_{kj}, b_{kj})$, $j = 1, \dots, q$ where the shape and scale parameters depend on the previous values of the chain. This proposal density is the full conditional of $\tilde{\sigma}_k$ for $\tilde{\phi} = 0$ and thus it leads to having higher acceptance rates for chain values with $\tilde{\phi}$ close to zero than in other cases.

Full conditional for ϕ

Using the same arguments as for σ we have

$$\pi(\phi | \alpha, \beta, \gamma, \sigma, \mathcal{C}, \mathcal{W}, \mathcal{Y}) \propto \prod_{k=1}^K \pi(\tilde{\phi}_k | \tilde{\beta}_k, \tilde{\sigma}_k, \mathcal{C}, \mathcal{W}, \mathcal{Y})$$

with

$$\pi(\tilde{\phi}_k | \tilde{\beta}_k, \tilde{\sigma}_k, \mathcal{C}, \mathcal{W}, \mathcal{Y}) \propto \pi(\tilde{\phi}_k) \prod_{i=1}^n \left\{ |\tilde{S}_{c_i}|^{-1/2} \exp \left[-\frac{1}{2} \sum_{t=1}^T \left(\tilde{u}_{it} - X'_t \tilde{\beta}_{c_i} \right)' S_{c_i, t}^{-1} \left(\tilde{u}_{it} - X'_t \tilde{\beta}_{c_i} \right) \right] \right\}^{I(c_i=k)},$$

Here $\pi(\tilde{\phi}_k)$ is the prior distribution for $\tilde{\phi}_k$. Barnard *et al.* (2000) proposed two alternative prior models for $\tilde{\phi}_k$. One is the marginally uniform prior, in which the marginal prior for each correlation is a modified beta distribution over $[-1, 1]$; with an appropriate choice of the beta parameters, this becomes a uniform marginal prior distribution. The other model for $\tilde{\phi}_k$ is called the jointly uniform prior. Liechty *et al.* (2004) discussed prior uniformly distributed over all possible correlation matrices. Under the bivariate case, $q = 2$, the choice is greatly simplified and we choose a uniform prior, i.e. $\phi_k \sim \mathcal{U}(-1, 1)$ for every k .

Even with this simple choice, the conditional simulation requires a M-H step. The independent proposal for ϕ_k is drawn from a uniform distribution “centered” around the current value in the chain, say ϕ_k^* , i.e.

$$\phi_k \sim U(\max\{-1, \phi_k^* - r\}, \min\{1, \phi_k^* + r\})$$

with $r > 0$. In the simulation experiments we have seen that for a value like as $r = 0.1$ the MCMC algorithm performs very well.

Full conditional for \mathcal{C}

We start by noting that

$$f(\mathcal{C}|\alpha, \boldsymbol{\beta}, \boldsymbol{\gamma}, \boldsymbol{\sigma}, \mathcal{W}, \mathcal{Y}) \propto \prod_{i=1}^n f(c_i|\alpha, \tilde{\beta}_{c_i}, \tilde{\sigma}_{c_i}, \tilde{\phi}_{c_i}, \mathbf{y}_i, \mathbf{w}_i)$$

The conditional probability of the membership of c_i is given by

$$f(c_i|\mathbf{y}_i, \mathbf{w}_i, \beta_{c_i}, \sigma_{c_i}, \alpha) = \frac{\alpha_{c_i} f(\mathbf{y}_i|\alpha, \tilde{\beta}_{c_i}, \tilde{\sigma}_{c_i}, \tilde{\phi}_{c_i}, \mathbf{w}_i, c_i)}{\sum_{k=1}^K \alpha_k f(\mathbf{y}_i|\alpha, \tilde{\beta}_k, \tilde{\sigma}_k, \tilde{\phi}_k, \mathbf{w}_i, k)} \quad (\text{S.3})$$

Full conditional for \mathcal{W}

We note that

$$\pi(\mathcal{W}|\alpha, \boldsymbol{\beta}, \boldsymbol{\gamma}, \boldsymbol{\sigma}, \boldsymbol{\phi}, \mathcal{C}, \mathcal{Y}) \propto \prod_{i=1}^n f(\mathbf{w}_i|\tilde{\beta}_{c_i}, \gamma_{c_i}, \tilde{\sigma}_{c_i}, \tilde{\phi}_{c_i}, c_i, \mathbf{y}_i) = \prod_{i=1}^n \prod_{t=1}^T f(\tilde{w}_{it}|\tilde{\beta}_{c_i}, \gamma_{c_i}, \tilde{\sigma}_{c_i}, \tilde{\phi}_{c_i}, c_i, \tilde{y}_{it})$$

with

$$f(\tilde{w}_{it}|\tilde{\beta}_{c_i}, \gamma_{c_i}, \tilde{\sigma}_{c_i}, \tilde{\phi}_{c_i}, c_i, \tilde{y}_{it}) \propto e(\tilde{w}_{it}; \gamma_{c_i},) |\tilde{S}_{c_i}|^{-1/2} \exp \left[-\frac{1}{2} \left(\tilde{u}_{it} - X'_t \tilde{\beta}_{c_i} \right)' S_{c_i, t}^{-1} \left(\tilde{u}_{it} - X'_t \tilde{\beta}_{c_i} \right) \right]$$

Here $e(\tilde{w}_{it}; \gamma)$ is a multivariate density function with marginal unit exponential distributions

and γ is a generic dependence parameter. The density $f(\tilde{w}_{it}|\tilde{\beta}_{c_i}, \gamma_{c_i}, \tilde{\sigma}_{c_i}, \tilde{\phi}_{c_i}, c_i, \tilde{y}_{it})$ seems not available in closed form for any reasonable choice of $e(\tilde{w}_{it}; \gamma)$. Once again we resort to a M-H step for simulating from it.

In case of independence of the components of \tilde{w}_{it} , conventionally identified with $\gamma_k = 0$, for all k , we have $e(\tilde{w}_{it}; 0) = \prod_{j=1}^q \exp(-w_{ijt})$. Moreover we pretend that $\tilde{\phi}_k = 0$ and in this case

$$f(\tilde{w}_{it}|\tilde{\beta}_{c_i}, 0, \tilde{\sigma}_{c_i}, 0, c_i, \tilde{y}_{it}) = \prod_{j=1}^q e(w_{ijt}|\tilde{\beta}_{c_i}, \sigma_{c_i,j}, c_i, y_{ijt})$$

where

$$\begin{aligned} e(\tilde{w}_{ijt}|\tilde{\beta}_{c_i}, \sigma_{c_i,j}, c_i, y_{ijt}) &\propto w_{ijt}^{-1/2} \exp \left\{ -\frac{1}{2\sigma_{c_i,j}^2 \omega^2 w_{ijt}} \left(y_{ijt} - X'_{jt} \tilde{\beta}_{c_i} - \theta w_{ijt} \right)^2 - \frac{w_{ijt}}{\sigma_{c_i,j}^2} \right\} \\ &\propto w_{ijt}^{-1/2} \exp \left[-\frac{1}{2\sigma_{c_i,j}^2 \omega^2 w_{ijt}} \left\{ (y_{it} - X'_{jt} \tilde{\beta}_{c_i})^2 + \theta^2 w_{ijt}^2 \right\} - \frac{w_{ijt}}{\sigma_{c_i,j}^2} \right] \\ &\propto w_{ijt}^{-1/2} \exp \left\{ -\frac{1}{2\sigma_{c_i,j}^2 \omega^2 w_{ijt}} (y_{it} - X'_{jt} \tilde{\beta}_{c_i})^2 - \frac{\theta^2 w_{ijt}}{2\sigma_{c_i,j}^2 \omega^2} - \frac{w_{ijt}}{\sigma_{c_i,j}^2} \right\} \\ &\propto w_{ijt}^{-1/2} \exp \left\{ -\frac{1}{2\sigma_{c_i,j}^2 \omega^2 w_{ijt}} (y_{it} - X'_{jt} \tilde{\beta}_{c_i})^2 - \frac{w_{ijt}}{2\sigma_{c_i,j}^2 \omega^2} (\theta^2 + 2\omega^2) \right\} \\ &\propto w_{ijt}^{-1/2} \exp \left\{ -\frac{1}{2\sigma_{c_i,j}^2 \omega^2 w_{ijt}} (y_{it} - X'_{jt} \tilde{\beta}_{c_i})^2 - \frac{\omega^2 w_{ijt}}{8\sigma_{c_i,j}^2} \right\} \\ &\propto w_{ijt}^{-1/2} \exp \left\{ -\frac{1}{2} \left(\frac{\omega^2 w_{ijt}}{4\sigma_{c_i,j}^2} + \frac{(y_{it} - X'_{jt} \tilde{\beta}_{c_i})^2}{\sigma_{c_i,j}^2 \omega^2 w_{ijt}} \right) \right\}. \end{aligned}$$

This expression resembles a Generalized Inverse Gaussian distribution $\text{GIG}(p, a, b)$ with the density

$$h(w; a, b, p) \propto w^{(p-1)} \exp \left\{ -\frac{1}{2} \left(aw + \frac{b}{w} \right) \right\}$$

where $a = \frac{\omega^2}{4\sigma_{c_i,j}^2}$, $b = \frac{(y_{ijt} - X'_{jt} \tilde{\beta}_{c_i})^2}{\sigma_{c_i,j}^2 \omega^2}$ and $p = 1/2$. Jorgensen (1982) notes that if $W \sim \text{GIG}(-p, b, a)$ then $1/W \sim \text{GIG}(p, a, b)$ and further that the distribution $\text{GIG}(-1/2, a, b)$

equals the Inverse Gaussian distribution

$$\begin{aligned} g(w|\mu, \lambda) &\propto w^{-3/2} \exp\left(\frac{-\lambda(w - \mu)^2}{2\mu^2 w}\right) \\ &\propto w^{-3/2} \exp\left(-\frac{\lambda w}{2\mu^2} - \frac{\lambda}{2w}\right) \end{aligned}$$

with $a = \lambda/\mu^2$, $b = \lambda$ and $p = -1/2$.

By recognizing these facts we can sample from the conditional posterior of w_{ijt} by drawing from the Inverse Gaussian distribution with $\lambda = \frac{\omega^2}{4\sigma_{c_i,j}^2}$ and $\frac{\lambda}{\mu^2} = \frac{(y_{ijt} - X'_{jt}\tilde{\beta}_{c_i})^2}{\sigma_{c_i,j}^2 \omega^2}$

that implies $\mu = \frac{\omega^2}{2|y_{ijt} - X'_{jt}\tilde{\beta}_{c_i}|}$.

Therefore, by choosing $e(\tilde{w}_{ijt}|\tilde{\beta}_{c_i}, \sigma_{c_i,j}, c_i, y_{ijt})$ as proposal density for every component w_{ijt} of \tilde{w}_{it} , we shall propose a value w such that $1/w$ is drawn from the density $\text{InvGauss}\left(\frac{\omega^2}{2|y_{ijt} - X'_{jt}\tilde{\beta}_{c_i}|}, \frac{\omega^2}{4\sigma_{c_i,j}^2}\right)$.

Full conditional for γ

We have

$$\pi(\gamma|\alpha, \beta, \sigma, \phi, \mathcal{C}, \mathcal{W}, \mathcal{Y}) \propto \prod_{k=1}^K \pi(\gamma_k|\mathcal{C}, \mathcal{W})$$

where

$$\pi(\gamma_k|\mathcal{C}, \mathcal{W}) \propto \pi(\gamma_k) \prod_{i=1}^n \left\{ \prod_{t=1}^T e(\tilde{w}_{it}; \gamma_{c_i}) \right\}^{I(c_i=k)}$$

Here $\pi(\gamma_k)$ is the prior distribution for γ_k . The density $\pi(\gamma_k|\mathcal{C}, \mathcal{W})$ seems not available in closed form for any reasonable choice of $e(\tilde{w}_{it}; \gamma)$. Once again we resort to a M-H step for simulating from it.