

Wavelet-Based Clustering for Mixed-Effects Functional Models in High Dimension

M. Giacomini,^{1,*} S. Lambert-Lacroix,^{2,**} G. Marot,^{3,4,5,***} and F. Picard^{6,****}

¹Laboratoire LJK, BP 53, Université de Grenoble et CNRS, 38041 Grenoble cedex 9, France

²UJF-Grenoble 1/CNRS/UPMF/TIMC-IMAG UMR 5525, Grenoble, F-38041, France

³Projet BAMBOO, INRIA Rhône-Alpes, F-38330 Montbonnot Saint-Martin, France

⁴Biostatistics, EA 2694, UDSL, Université Lille Nord de France

⁵MODAL, INRIA Lille Nord Europe, F-59650 Villeneuve d'Ascq, France

⁶LBBE, UMR CNRS 5558 Université Lyon 1, F-69622, Villeurbanne, France

* *email*: madison.giacomini@imag.fr

** *email*: sophie.lambert@imag.fr

*** *email*: guillemette.marot@univ-lille2.fr

**** *email*: franck.picard@univ-lyon1.fr

SUMMARY. We propose a method for high-dimensional curve clustering in the presence of interindividual variability. Curve clustering has long been studied especially using splines to account for functional random effects. However, splines are not appropriate when dealing with high-dimensional data and can not be used to model irregular curves such as peak-like data. Our method is based on a wavelet decomposition of the signal for both fixed and random effects. We propose an efficient dimension reduction step based on wavelet thresholding adapted to multiple curves and using an appropriate structure for the random effect variance, we ensure that both fixed and random effects lie in the same functional space even when dealing with irregular functions that belong to Besov spaces. In the wavelet domain our model resumes to a linear mixed-effects model that can be used for a model-based clustering algorithm and for which we develop an EM-algorithm for maximum likelihood estimation. The properties of the overall procedure are validated by an extensive simulation study. Then, we illustrate our method on mass spectrometry data and we propose an original application of functional data analysis on microarray comparative genomic hybridization (CGH) data. Our procedure is available through the R package `curvclust` which is the first publicly available package that performs curve clustering with random effects in the high dimensional framework (available on the CRAN).

KEY WORDS: Clustering; Functional data; Mixed models; Wavelets.

1. Introduction

Functional data analysis has gained increased attention in the past years, in particular in high-throughput biology with the use of mass spectrometry. This method is used to characterize the protein content of biological samples by separating compounds according to their mass to charge ratio (m/z). Among different technologies matrix assisted laser desorption and ionization, time-of-flight (MALDI-TOF) mass spectrometry is one the most used and has become standard to improve proteomic profiling of diseases as well as clinical diagnosis.

Dedicated methods have been developed to analyze such data for differential analysis, supervised classification and clustering (Hilario et al. 2006). Up to now the functional setting has mostly been developed for differential analysis (Morris et al. 2008). One central element is the modeling of the interindividual variability by using functional random effects, because subject-specific fluctuations are known to be the largest source of variability in mass-spec data (Eckel-Passow et al. 2009). In this article, we focus on the

nonsupervised task which consists in finding groups of individuals whose proteomic landscape is similar. Surprisingly the clustering task received less attention, and is mainly based on hierarchical clustering on the set of peaks detected across spectra (Bensmail et al. 2005; Morris et al. 2010). However, such method is known to depend heavily on the peak detection method and has the strong disadvantage to neglect the interindividual variability whereas this information should be central for subgroup discovery. Thus, our main focus in this article is modeling and clustering curves of this type in a functional mixed model framework.

When dealing with curve clustering in the presence of individual variability, a pioneer work is based on a spline decomposition of the signal (James and Sugar 2003) which resumes to a linear mixed effect model on which clustering and low-dimensional representation can be performed. However, splines show two main drawbacks: (i) they are inappropriate when dealing with functions that show peaks and irregularities, (ii) they require heavy computational efforts and so are not adapted to high dimensional data. On the contrary,

wavelet representations appear to be a natural framework to consider such irregularities through the sequence space of (usually sparse) Besov representation. Recent works have been done about estimation and inference in the functional mixed effects framework based on a wavelet decomposition approach. A fully Bayesian version has been proposed by Morris and Carroll (2006), with nonparametric estimates of fixed and random effects as well as between and within-curve covariance matrix estimates to accommodate a wide variety of correlation structures. In addition, Antoniadis and Sapatinas (2007) propose a study of both estimation and inference in a frequentist framework. In this article, we use a wavelet representation for both fixed and random effects to perform model-based clustering. Such strategy has been considered by Antoniadis, Bigot, and von Sachs (2008) and by Ray and Mallick (2006) without random effects for image clustering and for the analysis of time course experiments respectively. We use a similar approach and we extend it by adding functional random effects. Interindividual variability in the wavelet domain is modeled using results of Antoniadis and Sapatinas (2007) but accommodates a broader range of correlation structure. In particular we allow within curve correlation to vary over groups and positions. Then we propose a two-step procedure which involves a dimension reduction step and a clustering step based on the EM-algorithm. We also propose a model-selection criterion that accounts for the interindividual variability, and we define a rigorous simulation framework for curve clustering. Our method is implemented within the **R** package `curvclust` (available on the CRAN), which is the first available software dedicated to this task. In a first application, we illustrate our method on the mass spectrometry data first published in Petricoin et al. (2002).

Then our last contribution is to extend the use of functional models to another type of high throughput data which are comparative genomic hybridization (CGH) data. The CGH array technology is used to map copy number imbalances between genomes by hybridizing differentially labeled genomic DNAs on a chip. Fluorescence ratios are usually analyzed using change-point models to detect segments that correspond to homogeneous regions on the genome in terms of copy number. Clustering patients based on their CGH profiles is very promising and has been successfully used to identify molecular subtypes of cancer. However, clustering CGH profiles based on a segmentation has the same drawbacks that clustering mass spectra based on detected peaks: results depend on the segmentation methods. Moreover the interindividual variability has never been investigated in this type of data, whereas it is likely to represent an important part of the variability of the data especially for cancer profiles. We use the breast cancer data of Fridlyand et al. (2006) that have already been analyzed for nonsupervised clustering by Van Wieringen et al. (2008). We show the interest of functional random effects for these type of data and we discuss the impacts in terms of analysis and design for copy number studies.

2. Functional Clustering Modeling using Wavelets

2.1 Presentation of the Model

We observe N curves $Y_i(t)$ over M equally spaced time points $\mathbf{t} = (t_1, \dots, t_M)$ with $t_j \in [0, 1]$ for $j \in [1, M]$, and

$M = 2^J$ for some integer J . In the functional clustering setting we suppose that individuals are spread among L unknown clusters of prior size π_ℓ , $\ell = 1, \dots, L$, and we denote by $\zeta_{i\ell}$ the indicator variable that equals 1 if the i th individual is in the ℓ th group. Then, we consider the linear functional model such that given $\{\zeta_{i\ell} = 1\}$, $Y_i(t) = \mu_\ell(t) + E_i(t)$, where $\mu_\ell(t)$ is the principal functional fixed effect that characterizes cluster ℓ , $E_i(t)$ is a zero mean Gaussian process with covariance kernel $\text{cov}(E_i(t), E_i(t')) = \sigma_E^2 \delta_{tt'}$, where $\delta_{tt'}$ stands for the Kronecker product. In the following, we will use notations $\mathbf{Y}_i(\mathbf{t}) = (Y_i(t_1), \dots, Y_i(t_M))^T$, $\mu_\ell(\mathbf{t}) = (\mu_\ell(t_1), \dots, \mu_\ell(t_M))^T$ and $\mathbf{E}_i(\mathbf{t}) = (E_i(t_1), \dots, E_i(t_M))^T$. To handle subject-specific random deviations from the cluster average curve we introduce random functions $U_i(t)$ that are modeled as centered Gaussian processes with kernel $K_\ell(t, t') = \text{cov}(U_i(t), U_i(t'))$ (given $\{\zeta_{i\ell} = 1\}$), not necessarily stationary, but independent from $E_i(t)$. Then given $\{\zeta_{i\ell} = 1\}$, the previous model becomes $Y_i(t) = \mu_\ell(t) + U_i(t) + E_i(t)$ (2.1). Once defined in the functional domain, a classical approach is to convert the original infinite-dimensional clustering problem into a finite-dimensional problem using a functional basis representation of the model. At this step James and Sugar (2003) propose a spline-based representation of model (2.1) with individuals observed at sparse sets of time points like in longitudinal data. Our procedure is more adapted to high dimensional data thanks to the computational efficiency of wavelets, unlike splines that require matrix inversions whose complexity increases with the density of the design. Moreover, as we will see below, the wavelet representation allows us to account for a wider range of functional shapes than splines, thanks to their connection with Besov spaces. Using a wavelet representation of this model allows us to characterize different types of smoothness conditions assumed on the response curves $Y_i(t)$ by the mean of their wavelet coefficients. Moreover, wavelet representations are sparse for a wide variety of functional spaces, which is crucial when dealing with high dimensional data. This property will be central while performing dimension reduction. Briefly, we are working with a dyadic orthonormal wavelet basis $\{\phi_{j_0k}(t), k = 0, 1, \dots, 2^{j_0} - 1; \psi_{jk}(t), j \geq j_0, k = 0, \dots, 2^j - 1\}$ generated from a father wavelet ϕ and a mother wavelet ψ of regularity r , ($r \geq 0$). In this basis $Y_i(t)$ has the following decomposition: $Y_i(t) = \sum_{k=0}^{2^{j_0}-1} c_{i,j_0k}^* \phi_{j_0k}(t) + \sum_{j \geq j_0} \sum_{k=0}^{2^j-1} d_{i,jk}^* \psi_{jk}(t)$. In practice we use the discrete wavelet transform (DWT) which can be performed thanks to Mallat's fast algorithm with $\mathcal{O}(M)$ operations only. We denote by \mathbf{W} the $[M \times M]$ -matrix containing filters of the chosen wavelet basis. The resulting scaling and wavelet coefficients $\mathbf{c}_i = (c_{i,j_0k})_{k=0 \dots 2^{j_0}-1}$ and $\mathbf{d}_i = (d_{i,jk})_{j=j_0 \dots J-1}^{k=0 \dots 2^j-1}$ of the individual curves are empirical discrete coefficients. They are related to their theoretical continuous counterparts c_{i,j_0k}^* and $d_{i,jk}^*$ by: $c_{i,j_0k} \approx \sqrt{M} c_{i,j_0k}^*$ and $d_{i,jk} \approx \sqrt{M} d_{i,jk}^*$. In the following, we denote by $\alpha_\ell = (\alpha_{\ell,j_0k})_{k=0 \dots 2^{j_0}-1}$ and $\beta_\ell = (\beta_{\ell,jk})_{j=j_0 \dots J-1}^{k=0 \dots 2^j-1}$ the $[2^{j_0} \times 1]$ and $[(M - 2^{j_0}) \times 1]$ vectors of scaling and wavelet coefficients of $\mu_\ell(\mathbf{t})$, and we denote by $\nu_i = (\nu_{i,j_0k})_{k=0 \dots 2^{j_0}-1}$ and $\theta_i = (\theta_{i,jk})_{j=j_0 \dots J-1}^{k=0 \dots 2^j-1}$ the $[2^{j_0} \times 1]$ and $[(M - 2^{j_0}) \times 1]$ vectors of scaling and wavelet random coefficients of $\mathbf{U}_i(\mathbf{t}) = (U_i(t_1), \dots, U_i(t_M))^T$. We apply the DWT to model (2.1) such

that $\mathbf{WY}_i(\mathbf{t}) = \mathbf{W}\mu_\ell(\mathbf{t}) + \mathbf{WU}_i(\mathbf{t}) + \mathbf{WE}_i(\mathbf{t})$, and in the coefficients domain our model resumes to a linear mixed-effect model, such that given $\{\zeta_{i\ell} = 1\}$, $(\mathbf{c}_i^T, \mathbf{d}_i^T)^T = (\alpha_\ell^T, \beta_\ell^T)^T + (\nu_i^T, \theta_i^T)^T + (\varepsilon_{\mathbf{c}_i}^T, \varepsilon_{\mathbf{d}_i}^T)^T$. $(\varepsilon_{\mathbf{c}_i}^T, \varepsilon_{\mathbf{d}_i}^T)^T$ stands for the vector of errors on scaling and wavelet coefficients, distributed as $\mathcal{N}(0_M, \sigma_\varepsilon^2 \mathbf{I}_M)$ with 0_M the vector of zeros and \mathbf{I}_M the identity matrix of size M , and $\sigma_\varepsilon^2 = \sigma_E^2$. Then we suppose that $(\nu_i^T, \theta_i^T)^T \sim \mathcal{N}(0_M, \mathbf{G} = \text{Diag}(\mathbf{G}_\nu, \mathbf{G}_\theta))$, with, \mathbf{G}_ν and \mathbf{G}_θ the covariance matrices of ν_i and θ_i , respectively. We further suppose that these random coefficients are independent from the errors and that matrix \mathbf{G} is diagonal, thanks to the whitening property of wavelets (Zhang and Walter 1994). Without loss in generality, we will assume that $j_0 = 0$ in the following.

2.2 Besov Spaces and Specification of the Variance of Random Effects

The strength of the wavelet representation is that it allows us to handle very diverse shapes of curves among which curves with irregularities that lie in particular Besov spaces. Besov spaces consist of functions that have a specific degree of smoothness. Roughly speaking, for a Besov space $B_{p,q}^s[0,1]$, parameter s indicates the number of function's derivatives, where their existence is required in a L^p -sense, q allowing finer control of the function's regularity. For a detailed study of Besov spaces, we refer to Donoho and Johnstone (1998). When dealing with functional mixed models the difficulty lies in the control of the regularity of random functions U_i , so that if the fixed function μ_ℓ is supposed to belong to some Besov space, U_i belongs to the same functional space. Following Antoniadis and Sapatinas (2007), this goal is achieved by controlling the exponential decrease of the variances of the random wavelet coefficients such that $\mathbb{V}(\theta_{i,jk}) = 2^{-j\eta} \gamma_\theta^2$ with parameter η being associated with the regularity of process U_i . Indeed, Abramovich, Sapatinas, and Silverman (1998) state that given a mother wavelet ψ of regularity r , where $\max(0, \frac{1}{p} - \frac{1}{2}) < s < r$ and given that $\mu_\ell(t) \in B_{p,q}^s[0,1]$, then,

$$U_i(t) \in B_{p,q}^s[0,1] \text{ a.s.} \iff \begin{cases} s + \frac{1}{2} - \frac{\eta}{2} = 0 & \text{if } 1 \leq p < \infty \text{ and } q = \infty, \\ s + \frac{1}{2} - \frac{\eta}{2} < 0 & \text{otherwise.} \end{cases}$$

We further allow γ_θ^2 to depend on scale and position ($\gamma_{\theta,jk}^2$) as proposed by Morris and Carroll (2006) or on cluster ($\gamma_{\theta,\ell}^2$) or on both ($\gamma_{\theta,\ell,jk}^2$). As mentioned by Antoniadis and Sapatinas (2007), even if the model restricts matrix \mathbf{G} to the class of matrices diagonalisable by the DWT, modeling $\mathbb{V}(\theta_{i,jk})$ as a function of scale and position allows us to account for dependencies and nonstationarities in the functional domain.

2.3 Dimensionality Reduction

Wavelet representations are sparse for a wide class of functional spaces which makes their use very efficient when dealing with high dimensional data. In the case of a single curve, shrinkage estimation and hard thresholding have been developed by Donoho and Johnstone (1994). Both methods present the double advantage to reduce dimensionality and to ensure good reconstruction properties. In the framework of curve clustering, our goal is to reduce the dimensionality of the problem to handle heavy datasets and not to find the optimal reconstruction rule. With this in mind we follow the

strategy proposed by Antoniadis et al. (2008) and we propose a dimension reduction procedure that proceeds in two steps,

(1) We first perform individual denoising to keep coefficients which contain individual-specific information. This is done by applying nonlinear wavelet hard thresholding of coefficients \mathbf{d}_i via an universal threshold as described in Donoho and Johnstone (1994). For recall, it consists in setting to zero coefficients $d_{i,jk}$ whose absolute value are below the universal threshold $\sigma\sqrt{2\log M}$. A traditional way to estimate σ is to take the average of the N robust individual noise variance estimates defined by the median absolute deviation ($\widehat{\sigma}_{\text{MAD}}$) of empirical wavelet coefficients at the finest resolution level $J-1$ divided by 0.6745. In our setting this quantity provides a robust estimation of the variance level at the finest resolution level, i.e., $\mathbb{V}(d_{i,J-1,k}) = 2^{-(J-1)\eta} \gamma_\theta^2 + \sigma_\varepsilon^2$.

(2) In a second part, we take the union set of wavelet coefficients that survived thresholding. This has the advantage to remove wavelet coefficients that are zero for all individuals, and hence which are not informative regarding to the clustering goal.

As a first remark, we can point that a mixed-model specific thresholding rule could be applied by taking an estimate of the global variance of the observations which is given by $\mathbb{V}(d_{i,jk}) = 2^{-j\eta} \gamma_\theta^2 + \sigma_\varepsilon^2$. Such a level dependent thresholding would lead to greater variance estimate and hence to a greater dimensionality reduction. Nevertheless its estimation would require estimates of both parameters σ_ε^2 and γ_θ^2 . This can be easily done when the individual labels are known. Otherwise, this estimation is a difficult task when individual labels are unknown because it leads to estimate variance from samples with different and unknown means. Moreover, simulations showed that the difference was negligible (not shown). Finally note that we do not use the third reduction step proposed by Antoniadis et al. (2008) which is dedicated to image segmentation.

3. Parameter Estimation and Model Selection

3.1 An EM Algorithm for Maximum Likelihood Estimation

Once projected in the wavelet domain, the clustering model resumes to a standard clustering model with additional random effects whose variance is of particular form. Thus, parameters are estimated by maximum likelihood using the EM algorithm. Both label variables ζ and random effects (ν, θ) are unobserved and the complete data log-likelihood can be written such that $\log \mathcal{L}(\mathbf{c}, \mathbf{d}, \nu, \theta, \zeta; \pi, \alpha, \beta, \mathbf{G}, \sigma_\varepsilon^2) = \log \mathcal{L}(\mathbf{c}, \mathbf{d} | \nu, \theta, \zeta; \pi, \alpha, \beta, \sigma_\varepsilon^2) + \log \mathcal{L}(\nu, \theta | \zeta; \mathbf{G}) + \log \mathcal{L}(\zeta; \pi)$. This likelihood can be easily computed thanks to the properties of mixed linear models: $((\mathbf{c}_i^T, \mathbf{d}_i^T)^T | (\nu_i^T, \theta_i^T)^T, \{\zeta_{i\ell} = 1\}) \sim \mathcal{N}((\alpha_\ell^T + \nu_i^T, \beta_\ell^T + \theta_i^T)^T, \sigma_\varepsilon^2 \mathbf{I}_M)$. The E-step consists in replacing the unobserved variables by their conditional expectation. Hence, cluster labels predictors $\widehat{\zeta}_{i\ell}$ are up-dated using *posterior* probabilities $\tau_{i\ell}$ such that,

$$\widehat{\zeta}_{i\ell}^{[h+1]} = \tau_{i\ell}^{[h+1]} = \frac{\pi_\ell^{[h]} f(\mathbf{c}_i, \mathbf{d}_i; \alpha_\ell^{[h]}, \beta_\ell^{[h]}, \mathbf{G}^{[h]} + \sigma_\varepsilon^{2[h]} \mathbf{I}_M)}{\sum_p \pi_p^{[h]} f(\mathbf{c}_i, \mathbf{d}_i; \alpha_p^{[h]}, \beta_p^{[h]}, \mathbf{G}^{[h]} + \sigma_\varepsilon^{2[h]} \mathbf{I}_M)},$$

with $f(\cdot)$ the probability density function of the Gaussian distribution. Then, using notation $\widehat{\nu}_{i\ell} = \mathbb{E}(\nu_i | \mathbf{c}_i, \zeta_{i\ell} = 1) = (\widehat{\nu}_{i,j_0 k \ell})_{k=0, \dots, 2^j - 1}$ and $\widehat{\theta}_{i\ell} = \mathbb{E}(\theta_i | \mathbf{d}_i, \zeta_{i\ell} = 1) = (\widehat{\theta}_{i,j k \ell})_{j=j_0, \dots, J-1}^{k=0, \dots, 2^j - 1}$, we apply the Henderson's trick (Robinson

1991) to get the following updates of the best linear unbiased predictors (BLUPs) of random effects: $\widehat{\nu}_{i\ell}^{[h+1]} = (\mathbf{c}_i - \alpha_{i\ell}^{[h]}) / (1 + \lambda_{\nu}^{[h]})$, and $\widehat{\theta}_{i\ell}^{[h+1]} = (\mathbf{d}_i - \beta_{i\ell}^{[h]}) / (1 + 2^{j\eta} \lambda_{\theta}^{[h]})$, with $(\lambda_{\nu}, \lambda_{\theta}) = (\sigma_{\varepsilon}^2 / \gamma_{\nu}^2, \sigma_{\varepsilon}^2 / \gamma_{\theta}^2)$. As for the maximization part, it provides the estimators of the mean curve coefficients $\alpha_{i\ell}^{[h+1]} = \sum_{i=1}^n \widehat{\zeta}_{i\ell}^{[h+1]} (\mathbf{c}_i - \widehat{\nu}_{i\ell}^{[h+1]}) / \widehat{N}_{i\ell}^{[h+1]}$, and $\beta_{i\ell}^{[h+1]} = \sum_{i=1}^n \widehat{\zeta}_{i\ell}^{[h+1]} (\mathbf{d}_i - \widehat{\theta}_{i\ell}^{[h+1]}) / \widehat{N}_{i\ell}^{[h+1]}$, with $\widehat{N}_{i\ell}^{[h+1]} = \sum_i \widehat{\zeta}_{i\ell}^{[h+1]}$, and $\pi_{i\ell}^{[h+1]} = \widehat{N}_{i\ell}^{[h+1]} / N$. Moreover, the EM algorithm provides a ML estimator of the variances of the model (using $j_0 = 0$): $N(M - 1) \gamma_{\theta}^{2[h+1]} = \sum_{i,j,k\ell} 2^{j\eta} \widehat{\zeta}_{i\ell}^{[h+1]} (\widehat{\theta}_{i,jk\ell}^{2[h+1]} + \frac{\sigma_{\varepsilon}^{2[h]}}{1 + 2^{j\eta} \lambda_{\theta}^{[h]}})$, $N \gamma_{\nu}^{2[h+1]} = \sum_{i\ell} \widehat{\zeta}_{i\ell}^{[h+1]} (\widehat{\nu}_{i,00\ell}^{2[h+1]} + \frac{\sigma_{\varepsilon}^{2[h]}}{1 + \lambda_{\nu}^{[h]}})$, and $MN \sigma_{\varepsilon}^{2[h+1]} = \sum_{i\ell} \widehat{\zeta}_{i\ell}^{[h+1]} \{ \sum_{j,k} [(d_{i,jk} - \widehat{\beta}_{i,jk}^{[h+1]} - \widehat{\theta}_{i,jk\ell}^{[h+1]})^2 + \frac{\sigma_{\varepsilon}^{2[h]}}{1 + \lambda_{\theta}^{[h]} 2^{j\eta}}] + (c_{i,00} - \widehat{\alpha}_{i,00}^{[h+1]} - \widehat{\nu}_{i,00\ell}^{[h+1]})^2 + \frac{\sigma_{\varepsilon}^{2[h]}}{1 + \lambda_{\nu}^{[h]}} \}$. Note that we use SEM, a stochastic version of EM to avoid random initializations (Celeux and Diebolt 1986). Hard clustering can also be performed using the Maximum a posteriori (MAP) rule based on posterior probabilities ($\tau_{i\ell}$). As last point, we mention that η can be estimated by maximization of the likelihood using the golden search section algorithm (Kiefer 1953).

3.2 Choosing the Number of Clusters

We propose to choose the number of clusters using the framework of penalized likelihoods. In the following, we use notations $\mathbf{m}_L[\gamma^2]$, $\mathbf{m}_L[\gamma_{\ell}^2]$ for clustering models with L groups with constant and group-dependent variances, respectively. We first use the Bayesian Information Criterion and we select the dimension that maximizes

$$\text{BIC}(\mathbf{m}_L[\gamma^2]) = \log \mathcal{L}(\mathbf{c}, \mathbf{d}; \widehat{\pi}, \widehat{\alpha}, \widehat{\beta}, \widehat{\mathbf{G}}, \widehat{\sigma}_{\varepsilon}^2, \mathbf{m}_L[\gamma^2]) - \frac{|\mathbf{m}_L[\gamma^2]|}{2} \times \log(N).$$

This classical criterion is a penalized version of the observed-data log-likelihood where $|\mathbf{m}_L[\gamma^2]| = (M + 1)L + |\mathbf{G}|$ is the number of free parameters of a model with L clusters, the dimension of \mathbf{G} (denoted by $|\mathbf{G}|$ here) depending on the variance structure of the random effects. When considering mixed models, it is likely that the prediction of the random effects provides information regarding the number of clusters to select. To use information from hidden variables we propose to derive an integrated classification likelihood criterion in the spirit of Biernacki, Celeux, and Govaert (2000). The ICL criterion is based on the integrated likelihood of the complete data: $\log \mathcal{L}(\mathbf{c}, \mathbf{d}, \nu, \theta, \zeta | \mathbf{m}_L[\gamma_{\ell}^2]) = \log \mathcal{L}(\mathbf{c}, \mathbf{d} | \nu, \theta, \zeta, \mathbf{m}_L[\gamma_{\ell}^2]) + \log \mathcal{L}(\nu, \theta | \zeta, \mathbf{m}_L[\gamma_{\ell}^2]) + \log \mathcal{L}(\zeta | \mathbf{m}_L[\gamma_{\ell}^2])$. For the first term we use a BIC-like approximation such that $-2 \log \mathcal{L}(\mathbf{c}, \mathbf{d} | \nu, \theta, \zeta, \mathbf{m}_L[\gamma_{\ell}^2]) \simeq NM \log \text{RSS}(\mathbf{c}, \mathbf{d} | \nu, \theta) + (ML + 1) \times \log(N)$, with $\text{RSS}(\mathbf{c}, \mathbf{d} | \nu, \theta, \zeta)$ the residual sum of squares defined such that $\text{RSS}(\mathbf{c}, \mathbf{d} | \nu, \theta, \zeta) = \sum_{i\ell} \zeta_{i\ell} \|\mathbf{c}_i - \widehat{\alpha}_{i\ell} - \nu_{i\ell}\|^2 + \sum_{i\ell} \zeta_{i\ell} \|\mathbf{d}_i - \widehat{\beta}_{i\ell} - \theta_{i\ell}\|^2$. Then we derive the integrated log-likelihood of the random effects. We assume a noninformative Jeffrey prior for the variance parameters such that $g(\gamma_{\nu,\ell}^2 | \zeta, \mathbf{m}_L[\gamma_{\ell}^2]) \propto 1/\gamma_{\nu,\ell}^2$. Using notations $N_{\ell} = \sum_{i=1}^N \zeta_{i\ell}$

and $\text{RSS}_{\ell}(\nu, \zeta) = \sum_{i=1}^N \zeta_{i\ell} \nu_{i,00\ell}^2$, we get,

$$\begin{aligned} -2 \log \mathcal{L}(\nu | \zeta, \mathbf{m}_L[\gamma_{\ell}^2]) &\simeq \sum_{\ell} N_{\ell} \log \text{RSS}_{\ell}(\nu, \zeta) \\ &- 2 \sum_{\ell} \log \Gamma(N_{\ell}/2). \end{aligned}$$

Similarly for the detail coefficients we get,

$$\begin{aligned} -2 \log \mathcal{L}(\theta | \zeta, \mathbf{m}_L[\gamma_{\ell}^2]) &\simeq (M - 1) \sum_{\ell} N_{\ell} \log \text{RSS}_{\ell}(\theta, \zeta) \\ &- 2 \sum_{\ell} \log \Gamma(N_{\ell}(M - 1)/2). \end{aligned}$$

Finally for the classification term a Dirichlet prior is assumed for π and the corresponding integrated likelihood is approximated such as,

$$\log \mathcal{L}(\zeta | \mathbf{m}_L[\gamma_{\ell}^2]) \simeq \sum_{\ell=1}^L N_{\ell} \log \left(\frac{N_{\ell}}{N} \right) - \frac{(L - 1)}{2} \log(N).$$

The last step of this derivation is to replace hidden variables by their predictions provided by the EM algorithm. Random effects (ν, θ) are replaced by their BLUP $(\widehat{\nu}, \widehat{\theta})$, and label variables ζ are replaced by their conditional expectation τ . Put together we obtain the following integrated classification likelihood criterion (ICL), such that $-2 \times \text{ICL}(\mathbf{m}_L[\gamma_{\ell}^2]) / N$ equals

$$\begin{aligned} M \log \text{RSS}(\mathbf{c}, \mathbf{d} | \widehat{\nu}, \widehat{\theta}, \tau) &+ \sum_{\ell} \widehat{\pi}_{\ell} \left[\log \text{RSS}_{\ell}(\widehat{\nu}, \tau) \right. \\ &+ (M - 1) \log \text{RSS}_{\ell}(\widehat{\theta}, \tau) \left. \right] \\ &- \frac{2}{N} \sum_{\ell} \left[\log \Gamma \left(\frac{\widehat{N}_{\ell}}{2} \right) + \log \Gamma \left(\frac{\widehat{N}_{\ell}(M - 1)}{2} \right) \right] \\ &- 2 \sum_{\ell=1}^L \widehat{\pi}_{\ell} \log(\widehat{\pi}_{\ell}) + \frac{(M + 1)L}{N} \times \log(N). \end{aligned}$$

Those criteria will be compared in the simulation study.

4. Simulations and Comparison of Methods

4.1 Definition of a General Simulation Framework

In this section, we propose to define a unified framework for synthetic data generation for functional mixed models and functional clustering models (FCMs). Using this unified strategy different methods can be fairly compared based on appropriately simulated data. First we properly define the signal-to-noise ratio (SNR) in the functional domain. The SNR is defined as the ratio of signal power to the power of the measurement noise corrupting the signal. In our case, the power of the signal is defined such as,

$$\begin{aligned} \lim_{T \rightarrow \infty} \frac{1}{T} \int_{\frac{T}{2}}^{-\frac{T}{2}} \sum_{\ell} \pi_{\ell} \mathbb{E}(|\mu_{\ell}(t) + U_i(t)|)^2 dt \\ = \frac{1}{M} \sum_{\ell=1}^L \pi_{\ell} \left(\sum_{k=0}^{2^{j_0}-1} \alpha_{\ell,j_0k}^2 + \sum_{j \geq j_0} \sum_{k=0}^{2^j-1} \beta_{\ell,jk}^2 \right) \\ + 2^{j_0} \gamma_{\nu}^2 + \frac{2^{j_0(1-\eta)} \gamma_{\theta}^2}{1 - 2^{(1-\eta)}}. \end{aligned}$$

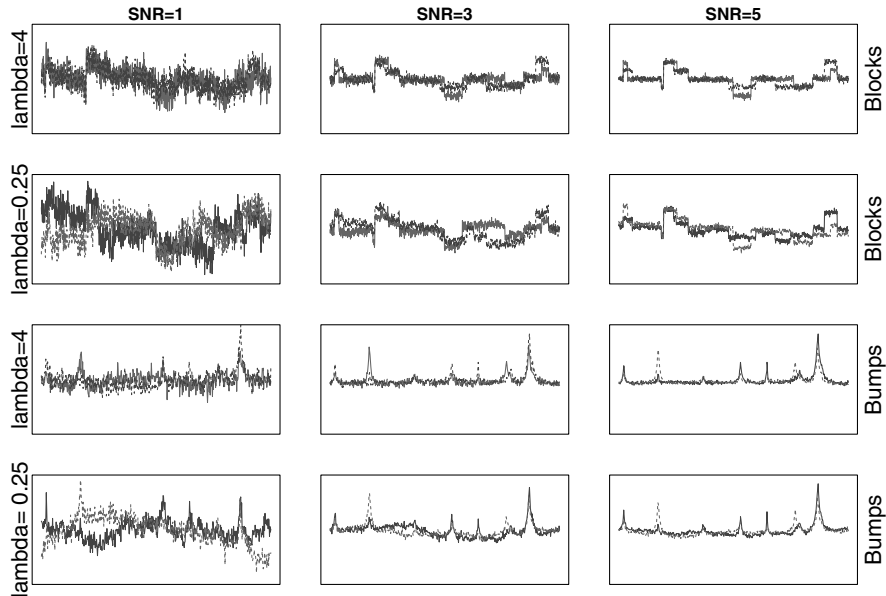


Figure 1. Example of simulated curves with varying SNR_μ and λ_U (one curve per cluster).

The derivation of such formula is given in the Web Supplementary Material. Hence we need to control two terms: SNR_μ that accounts for the power of the fixed effects and λ_U for the power of the random effect using an analogy with the λ parameter used in the EM algorithm. For this purpose we introduce parameters,

$$\text{SNR}_\mu^2 = \frac{1}{M\sigma_E^2} \sum_{\ell=1}^L \pi_\ell \left(\sum_{k=0}^{2^{j_0}-1} \alpha_{\ell,j_0k}^2 + \sum_{j \geq j_0} \sum_{k=0}^{2^j-1} \beta_{\ell,jk}^2 \right),$$

$$\lambda_U = \sigma_E^2 / \left(\gamma_\nu^2 + \frac{\gamma_\theta^2}{1 - 2^{-(1-\eta)}} \right).$$

When performing simulations, SNR_μ usually lies in $\{0.1, 1, 3, 5, 7\}$ and λ_U varies in $\{1/4, 1, 4\}$ such that small values of λ_U indicate an important variance for the random effects. In practice, we also choose $\gamma_\nu^2 = \gamma_\theta^2$.

To build fixed effects for simulations we generalize the approach described in Amato and Sapatinas (2005) which uses the well-known synthetic functions **Blocks**, **Bumps**, **Heavisine** and **Doppler** originally proposed by Donoho and Johnstone (1994). We choose L fixed effects for each synthetic function classes using expressions given in the Supplementary Material. Once parameters $(\text{SNR}_\mu, \lambda_U, \{\mu_\ell(t)\}_\ell)$ have been chosen (i.e., values for $\sigma_E^2, \gamma_\nu^2, \gamma_\theta^2$, and α_ℓ, β_ℓ are deduced), our simulation procedure is performed in the wavelet domain such that realizations of centered Gaussian distribution with variance $2^{-j\eta} \gamma_\theta^2$ are added to the fixed effect empirical wavelet coefficients to account for interindividual variability. Then Gaussian noise with variance σ_ε^2 is added to account for measurement errors. This unified method ensures that both fixed and random effects lie in the same Besov space, as mentioned earlier, and observed signals $\mathbf{Y}_i(\mathbf{t})$ can be recovered using the inverse DWT. An example of such simulated data is given in Figure 1.

4.2 Simulation Design and Indicators of Performance

Because too many configurations could be explored using simulations, we propose to fix the number of individuals at $N = 50$, the number of groups at $L = 2, 4$, the length of the signals at $M = 512$, and parameter η is set to 2. Then the simulation design explores the following configurations: $\text{SNR}_\mu \in \{0.1, 1, 3, 5, 7\}$, $\lambda_U \in \{1/4, 1, 4\}$, $\pi \in \{0.1, 0.25, 0.5\}$ ($\pi = 1/4$ when $L = 4$), each simulation being repeated 50 times. In terms of methods, we compare functional clustering models with or without mixed effects (FCMM/FCM, Functional Clustering Mixed Model/Functional Clustering Model), and we consider (or not) the dimension reduction method based on the union of coefficients. We compare these four methods to the functional clustering mixed model based on splines as proposed by James and Sugar (2003) whose R code is available on the web page of the authors (<http://www-bcf.usc.edu/gareth/>). Our purpose is to highlight the benefit of using wavelets when dealing with high dimensional data.

The performance of the clustering procedures are compared using the empirical error rate (EER) defined by $\text{EER} = 1/N \sum_{i=1}^N \sum_{\ell=1}^L \mathbb{I}\{\widehat{\zeta}_{i\ell}^{\text{MAP}} \neq \zeta_{i\ell}\}$, where $\widehat{\zeta}_{i\ell}^{\text{MAP}}$ is the predicted class for individual i using the MAP rule, and $\zeta_{i\ell}$ is the true class. This criteria ranges from 0, for which no classification error is made to 1 which means that all individuals are misclassified. We finally consider the speed of execution of each procedure.

4.3 Simulation Results

4.3.1 Clustering results. Figure 2 presents the variations of the Empirical Error Rates according to SNR_μ and to the strength of the random effect (a small λ_U indicates a strong random effect). A general comment is that the functional clustering mixed model (FCMM) outperforms all methods in terms of EER compared with the FCM and Splines. This result is true even for unbalanced clusters and with an increasing number of groups (see Supplementary Material). FCMM

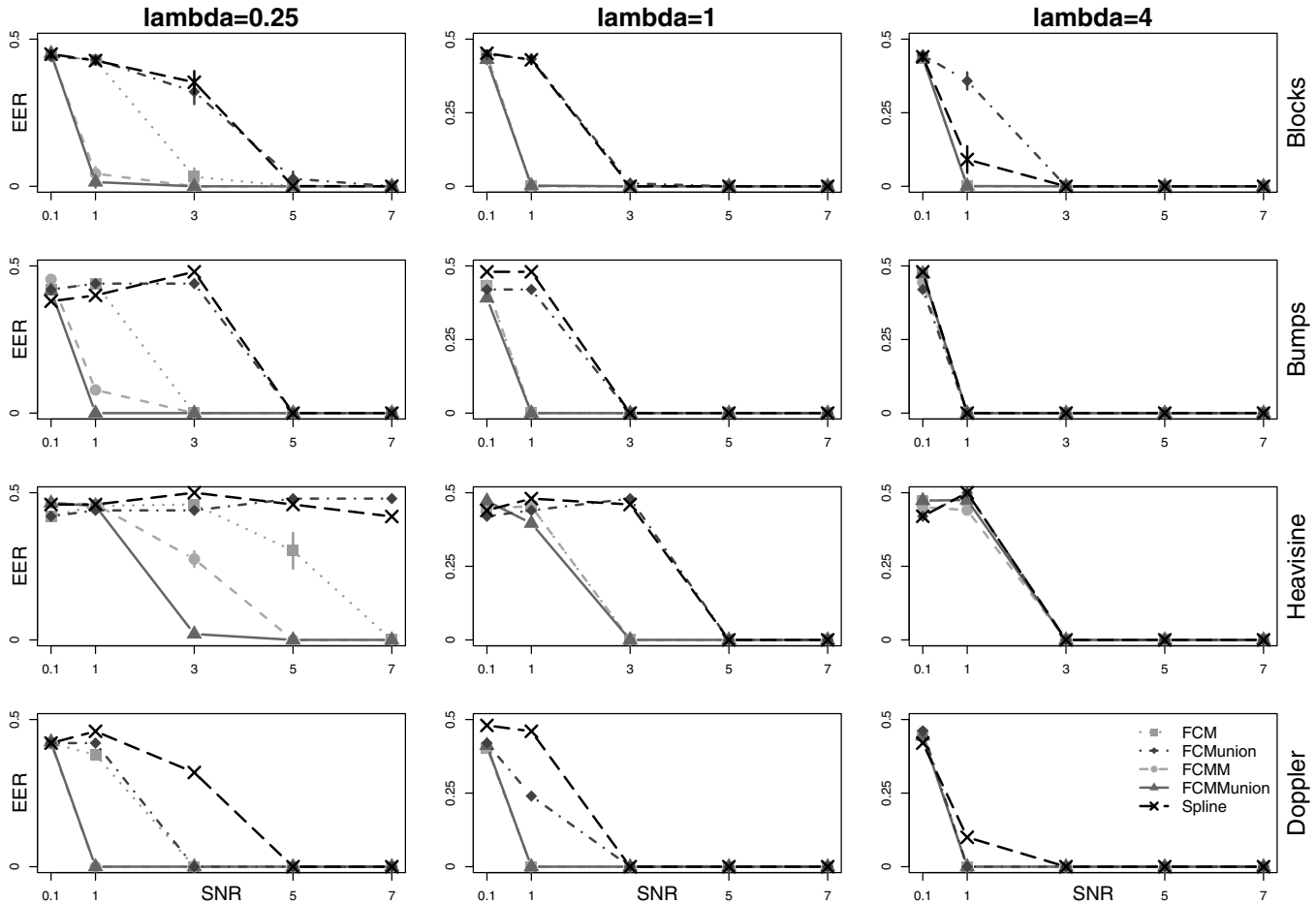


Figure 2. Variation of the empirical error rate (EER) for different estimation methods: functional clustering mixed model (FCMM), functional clustering model (FCM), with or without dimension reduction (“union”), and Splines. In columns different intensities for the variance of the random effect are considered: $\lambda_U = 0.25/1/4$ for a strong/mild/small random effect. In rows are considered different shapes for the mean curve of each group (Blocks, Bumps, Heavisine, Doppler). Results correspond to $L = 2$ clusters with balanced proportions 0.5/0.5

has two main advantages. First the modeling of functional random effects leads to a better identification of the informative structures in terms of clustering. Table 1 clearly shows that FCMM is the best method to estimate the variance of the residuals contrary to FCM that provides over-estimates (which leads to poor clustering performance).

Then dimension reduction increases the performance of FCMM by removing coefficients that are not informative with respect to clustering. This is not true for the FCM for which dimension reduction increases the EER. This trend can be explained by the bad estimation of the error’s variance when random effects are not considered in the model. The selection of the coefficients that all survived thresholding leads to worst estimators in the case of FCM but the impact is moderate on the FCMM (Table 1). In the Supplementary Material we also illustrate the performance of the dimension reduction procedure. This table was not provided by Antoniadis et al. (2008) when they first proposed the union-set method. Our results show that taking the union of coefficients that survived thresholding keeps less than 10% of the coeffi-

cients. Among those coefficients, we show that a high proportion should have been thresholded whereas they are not. This means that the procedure is sensitive but not very specific, as expected when considering a union-based strategy. However, because our objective is not functional reconstruction, we consider that keeping too many coefficients is not a major issue.

Our last point concerns the time of execution of each method. When dealing with high dimensional data, it is crucial to propose methods that show reasonable computational time. Table 1 clearly shows that using wavelet-based FCMMs gives the best execution times, and even when random effects are considered, time of execution remains moderate (less than 10 minutes for $N = 50$ individuals and $M = 512$ positions). Splines are known to be poorly efficient in terms of computational efficiency. This issue becomes critical when dealing with functional models with many individuals. The size of our simulated datasets was the upper limit that could be analyzed by Splines, in particular due to memory constraints. To this extent, our R package `curvclust` is the only freely

Table 1

Relative bias of the estimator of the error variance: $(\sigma^2 - \hat{\sigma}^2)/\sigma^2$, and average time of execution (TOE) in minutes for different models on simulated data ($N = 50$ individuals, $M = 512$ positions). FCM, functional clustering model, FCMM, functional clustering mixed model. FCMu/FCMMu, functional clustering (mixed) models based on the union of coefficients for dimension reduction. Programs were run on a cluster of 2 octo-bicore Opteron 2.8 Ghz and 2 octo-quadcore Opteron 2.3 GHz

SNR $_{\mu}^2$		Bias					TOE				
		0.1	1	3	5	7	0.1	1	3	5	7
FCM	Blocks	-2.57	-2.66	-2.96	-3.02	-2.99	2.3	2.4	2.3	2.4	2.3
	Bumps	-2.50	-2.69	-2.93	-2.93	-2.93	2.6	2.5	2.6	2.5	2.5
	Heavisine	-2.15	-2.17	-3.22	-4.30	-2.50	2.8	2.7	2.7	2.7	2.8
	Doppler	-2.73	-3.07	-3.32	-3.33	-3.33	2.9	3.2	3.1	3.2	3.2
FCMu	Blocks	-12.93	-11.33	-9.42	-9.38	-8.89	0.4	0.4	0.5	0.5	0.5
	Bumps	-12.98	-11.11	-13.46	-11.98	-11.93	0.5	0.5	0.5	0.5	0.5
	Heavisine	-11.62	-10.20	-10.07	-12.05	-15.68	0.5	0.5	0.5	0.5	0.5
	Doppler	-14.75	-13.14	-11.33	-8.59	-7.87	0.5	0.5	0.5	0.6	0.6
FCMM	Blocks	0.11	0.05	-0.01	-0.01	-0.00	16.0	16.1	15.6	15.8	16.0
	Bumps	0.09	0.04	0.01	0.01	0.01	16.1	16.3	15.2	15.3	15.4
	Heavisine	0.10	0.09	0.08	0.03	0.02	16.4	16.2	16.0	16.4	15.9
	Doppler	0.08	0.01	-0.02	-0.02	-0.01	17.5	17.4	17.5	16.4	17.0
FCMMu	Blocks	-0.11	-0.06	0.03	0.06	0.05	6.9	7.1	7.6	7.6	7.6
	Bumps	-0.10	-0.04	-0.08	-0.08	-0.05	6.7	6.7	6.8	6.7	6.7
	Heavisine	-0.10	-0.10	-0.18	-0.21	-0.19	7.1	7.3	6.8	6.8	6.8
	Doppler	-0.18	-0.06	-0.04	-0.16	-0.11	7.3	7.1	7.3	7.8	7.9
Spline	Blocks	25.5	26.2	23.0	23.6	22.3
	Bumps	23.3	26.6	22.0	21.2	21.7
	Heavisine	24.2	21.6	21.8	22.4	22.3
	Doppler	33.2	32.4	24.2	24.8	24.2

available software that performs curve clustering with functional random effects within a reduced amount of time in high dimension.

4.3.2 Model selection results. The model selection criteria are compared using the same simulation design with four groups (Figure 3). The BIC selects four clusters even when the SNR is low (except for Heavisine), contrary to ICL which is more stringent. Their behavior differ slightly with respect to the strength of the random effect, with ICL penalizing more when the random effect is strong whereas BIC gives similar results with respect to the strength of the random effect. Overall, differences between criteria are mild.

5. Applications

5.1 Mass Spectrometry Data

We first consider a SELDI-TOF mass spectrometry dataset issued from a study on ovarian cancer (Petricoin et al. 2002). The sample set includes serum profiles of 162 subjects with ovarian cancer and 91 non-cancer control subjects. Each serum profile consists of 15,154 recorded intensities corresponding to distinct m/z values. This dataset was produced by the Ciphergen WCX2 protein chip. It is available through the Clinical Proteomics Programs Databank (<http://home.ccr.cancer.gov/ncifdaproteomics/ppatterns.asp>, ovarian dataset 8-7-02). Before clustering, raw data are background corrected using a quantile regression procedure, and spectra are aligned using a procedure based on wavelets zero crossings (Antoniadis et al. 2007). Then the ovarian cancer dataset is made of 8192 intensities

within the range of m/z ratio [1500, 14,000], ratios below 1500 being discarded due to the effects of matrix. We compare wavelet-based FCMs on these data considering different random effect structures. Procedures are applied in a non-supervised framework to retrieve the known labels (cancer/control) and comparisons are based on empirical error rate estimates (EER, Table 2). Note that the spline-based procedure of James and Sugar (2003) could not be applied on these data because of their too high dimensionality.

The first result is that empirical error rates are high for all methods and that the introduction of random effects slightly decreases the EER whatever the random effect structure (from 38% to $\sim 25\%$). To investigate the origins of such modest performance, we also performed clustering based on group-wise aligned spectra instead of global alignment (which should be done in the unsupervised context). Results are striking: when spectra are aligned according to known labels, model $\mathbf{m}_2[\gamma_{jk}^2]$ (for which the variance of random effects depends on scale and position) results in one mismatch only (EER=0.4%). This result leads to the following conclusions. First spectra alignment is a challenge when performing subgroup discovery, and the task is much more difficult compared with supervised clustering for which labels are known. Indeed inaccuracy in spectra alignment could lead to artificial differences in individual serum profiles which decreases the performance of clustering. A promising (but challenging) perspective would be to perform clustering and alignment simultaneously. Moreover as wavelets have been shown to perform best for peak-detection/alignment (Yang et al. 2009),

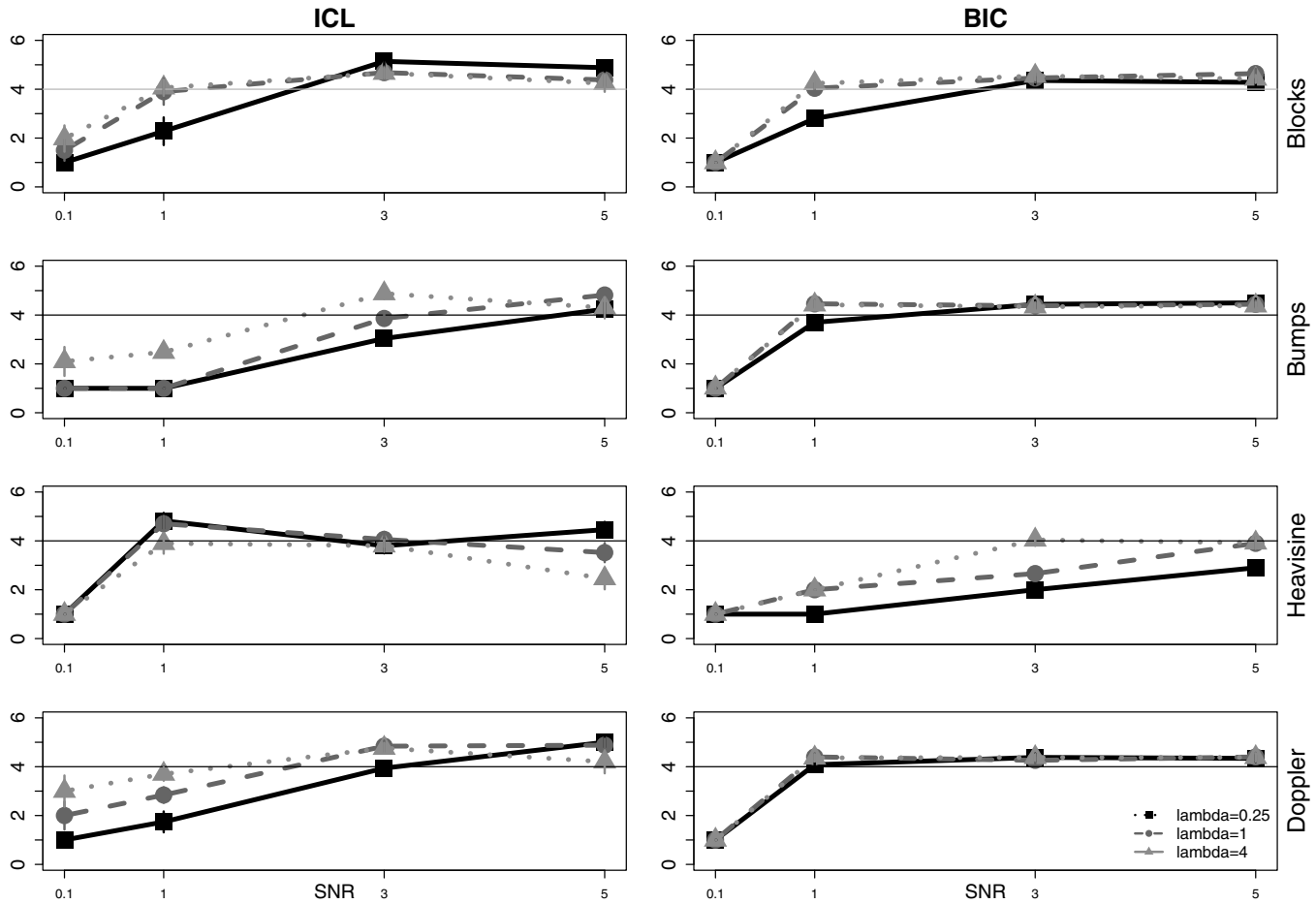


Figure 3. Estimated number of clusters using ICL and BIC when the simulated number of clusters is four (with balanced cluster sizes).

Table 2

Empirical error rates (in percent) for the Petricoin et al. (2002) data for different models: functional clustering without random effects, two groups (\mathbf{m}_2), functional clustering with random effect with different variance structures for the random effect: constant $\mathbf{m}_2[\gamma^2]$, group $\mathbf{m}_2[\gamma_\ell^2]$, scale-position $\mathbf{m}_2[\gamma_{jk}^2]$, or group-scale-position dependent $\mathbf{m}_2[\gamma_{\ell,jk}^2]$

	\mathbf{m}_2	$\mathbf{m}_2[\gamma^2]$	$\mathbf{m}_2[\gamma_\ell^2]$	$\mathbf{m}_2[\gamma_{jk}^2]$	$\mathbf{m}_2[\gamma_{\ell,jk}^2]$
Global alignment	38	24	24	23	23
Group alignment	20	21	22	0.4	36

our wavelet-based procedure for clustering would be a good starting point to integrate both strategies.

Then a second result is that best clustering performance are provided by a functional clustering mixed model for which the random effect has a covariance structure that depends on both scale and location. This implies that interindividual variations occur at specific ranges of m/z values, which reinforces the importance of correct spectra alignment. Interestingly, an important proportion of variance terms are close to zeros which would make the BLUPs sparse if dimension reduction was per-

formed on random effects. Unfortunately, the task is difficult in the nonsupervised setting because BLUPs can not be computed without the knowledge of group-specific means (which would be possible in the supervised setting). Thus dimension reduction for clustering using mixed functional model remains challenging and still needs to be investigated.

5.2 Comparative Genomic Hybridization Data

In this last application we consider the clustering of breast-cancer tumors based on their copy number aberration profiles measured by array-based Comparative Genomic Hybridization (Fridlyand et al. 2006). Array CGH is a widely used technology that enables the characterization of genome-wide chromosomal aberrations using the microarray technology. Many statistical methods have been developed to analyze these data (van de Wiel et al. 2011). They are mainly based on segmentation methods to retrieve segments of homogeneous copy number along the genome.

Clustering individuals based on their CGH profiles is a very challenging issue and has already been considered to identify new subtypes of tumors (Chin et al. 2007). For now, subgroup discovery is mainly performed using hierarchical clustering based on segmentation results (Van Wieringen et al. 2008). However, the interindividual variability has never been

Table 3
*Estimated SNR_μ^2 and λ_U for the breast tumor dataset of
 Fridlyand et al. (2006)*

Cluster ID	Complete dataset	
	$\widehat{\text{SNR}}_\mu^2$	$\widehat{\lambda}_U$
1	2.1e-4	3.9e-04
2	2.3e-3	3.8e-05
3	1.3e-3	6.4e-04
4 (1q/16p)	1.5e-3	1.3e-04
5	9.3e-4	4.3e-05
ER+ dataset		
1	2.1e-3	2.2e-04
2	7.8e-3	1.9e-05
3	1.1e-2	3.8e-05
4 (1q/16p)	4.4e-3	4.4e-04

quantified in these data, contrary to mass spectrometry for instance. Thus using our method for clustering with the Haar basis (piece-wise constant basis) is a way to perform subgroup discovery by considering random effects. In the Fridlyand et al. (2006) article, the authors analyzed the genomic profiles of 62 samples using P1/BAC CGH arrays (2464 genomic clones). We used the 55 profiles for which additional clinical information were available (the raw data can be downloaded as a supplementary material of the Fridlyand et al. 2006 article). The authors identified three main subtypes of breast cancer that differ with respect to level of genomic instability. Interestingly, Van Wieringen et al. (2008) re-analyzed the data and do not mention much correspondance between the two clustering results. Moreover, they discovered much more subgroups and noticed that “the samples in the study could be more heterogeneous than previously implied.”

We also find more subgroups than the original study, with five clusters selected by ICL (two by the BIC). First, this shows the power which is gained when considering the random effect in the selection step. Then we were able to identify the 1q/16p subtype on the complete dataset (with one mismatch). This subtype was identified in the first study (Fridlyand et al. 2006) but not by other clustering methods (Van Wieringen et al. 2008) whereas it is associated to the best patient outcome. Because two of the three identified clusters in the original article concern ER positive tumors, we also performed our method on this subset of patients and retrieve the 1q/16p subtype without mismatch. In this classification, one cluster was made of three tumors (S0041, S0041, S1519) also identified as similar in the original article. As a last result Table 3 indicates that the estimated signal to noise ratio is low and the impressive strength of the random effect ($\widehat{\lambda}_U \sim 10^{-4}$) also indicates that the interindividual variability is ultra-high in these data. As a consequence, finding clusters with biological significance will require rather hundreds/thousands of patients compared with 55 in the original study.

6. Conclusion

In this work we provide a methodology for model-based clustering of functional data in the presence of interindividual

variability. Our method is based on a wavelet decomposition of the signal and on a mixture model that integrates random effects. We illustrate the power of such an approach in two different fields of high-throughput biology using our package `curvclust`, and we show the potentialities of functional models on array CGH data. Overall, random effects allow us to properly model the variance structure of the data, and to exhibit the high proportion of variance due to interindividual variability. This part is usually omitted in high-throughput modelling. First perspective will concern the generalization of our approach to the supervised setting. Finding biomarkers has received enormous attention in the past years, with moderate success due to the lack of reproducibility. Our study in the nonsupervised framework shows that the interindividual variability is important in these data, which may be one explanation of the difficulty to find reliable markers. Integrating random effects in the supervised setting may produce more moderate results, but at least they would be more representative of the biological variability. Finally methodological perspectives of this work will mainly concern dimension reduction. The task is difficult in the non-supervised setting and the illustration on MS data shows that dimension reduction should be performed for fixed *and* for random effects which remains challenging. This would provide a better representation of the signal by thresholding coefficients with poor information, and would increase the speed of the estimation algorithm that is sensitive to the number of selected coefficients, which is of central interest of high dimensional data.

7. Supplementary Materials

Web Appendices, tables and figures referenced in Sections 4.1, 4.2 and 4.3 are available with this article at the Biometrics website on Wiley Online Library.

ACKNOWLEDGEMENTS

Part of this work was supported by the Interuniversity Attraction Pole (IAP) research network in Statistics P5/24 and by the MSTIC project of the Joseph-Fourier University.

REFERENCES

- Abramovich, F., Sapatinas, T., and Silverman, B. (1998). Wavelet thresholding via a Bayesian approach. *Journal of the Royal Statistical Society Series B Statistical Methodology* **60**, 725–749.
- Amato, U. and Sapatinas, T. (2005). Wavelet shrinkage approaches to baseline signal estimation from repeated noisy measurements. *Advances and Applications in Statistics* **51**, 21–50.
- Antoniadis, A., Bigot, J., Lambert-Lacroix, S., and Letue, F. (2007). Non parametric pre-processing methods and inference tools for analyzing time-of-flight mass spectrometry data. *Current Analytical Chemistry* **3**, 127–147.
- Antoniadis, A., Bigot, J., and von Sachs, R. (2008). A multiscale approach for statistical characterization of functional images. *Journal of Computational and Graphical Statistics* **18**, 216–237.
- Antoniadis, A. and Sapatinas, T. (2007). Estimation and inference in functional mixed-effects models. *Computational Statistics & Data Analysis* **51**, 4793–4813.
- Bensmail, H., Aruna, B., Semmes, O. J., and Haoudi, A. (2005). Functional clustering algorithm for high-dimensional proteomics data. *Journal of Biomedicine and Biotechnology* **2005**, 80–86.

- Biernacki, C., Celeux, G., and Govaert, G. (2000). Assessing a mixture model for clustering with the integrated completed likelihood. *IEEE PAMI* **22**, 719–725.
- Celeux, G. and Diebolt, J. (1986). The SEM algorithm: A probabilistic teacher algorithm derived from the EM algorithm for the mixture problem. *Computational Statistics Quarterly* **2**, 73–82.
- Chin, S. F., Teschendorff, A. E., Marioni, J. C., Wang, Y., Barbosa-Morais, N. L., Thorne, N. P., Costa, J. L., Pinder, S. E., van de Wiel, M. A., Green, A. R., Ellis, I. O., Porter, P. L., Tavaré, S., Brenton, J. D., Ylstra, B., and Caldas, C. (2007). High-resolution aCGH and expression profiling identifies a novel genomic subtype of ER negative breast cancer. *Genome Biology* **8**, R215.
- Donoho, D. and Johnstone, I. (1994). Ideal spatial adaptation by wavelet shrinkage. *Biometrika* **81**, 425–455.
- Donoho, D. and Johnstone, I. (1998). Minimax estimation via wavelet shrinkage. *Annals of Statistics* **26**, 879–921.
- Eckel-Passow, J. E., Oberg, A. L., Therneau, T. M., and Bergen, H. R. (2009). An insight into high-resolution mass-spectrometry data. *Biostatistics* **10**, 481–500.
- Fridlyand, J., Snijders, A. M., Ylstra, B., Li, H., Olshen, A., Segraves, R., Dairkee, S., Tokuyasu, T., Ljung, B. M., Jain, A. N., McLennan, J., Ziegler, J., Chin, K., Devries, S., Feiler, H., Gray, J. W., Waldman, F., Pinkel, D., and Albertson, D. G. (2006). Breast tumor copy number aberration phenotypes and genomic instability. *BMC Cancer* **6**, 96.
- Hilario, M., Kalousis, A., Pellegrini, C., and Muller, M. (2006). Processing and classification of protein mass spectra. *Mass Spectrometry Reviews* **25**, 409–449.
- James, G. and Sugar, C. (2003). Clustering for sparsely sampled functional data. *Journal of the American Statistical Association* **98**, 397–408.
- Kiefer, J. (1953). Sequential minimax search for a maximum. *Proceedings of the American Mathematical Society* **4**, 502–506.
- Morris, J. S., Baggerly, K. A., Gutstein, H. B., and Coombes, K. R. (2010). Statistical contributions to proteomic research. *Methods in Molecular Biology* **641**, 143–166.
- Morris, J. S., Brown, P. J., Herrick, R. C., Baggerly, K. A., and Coombes, K. R. (2008). Bayesian analysis of mass spectrometry proteomic data using wavelet-based functional mixed models. *Biometrics* **64**, 479–489.
- Morris, J. S. and Carroll, R. J. (2006). Wavelet-based functional mixed models. *Journal of the Royal Statistical Society Series B Statistical Methodology* **68**, 179–199.
- Petricoin, E. F., Ardekani, A. M., Hitt, B. A., Levine, P. J., Fusaro, V. A., Steinberg, S. M., Mills, G. B., Simone, C., Fishman, D. A., Kohn, E. C., and Liotta, L. A. (2002). Use of proteomic patterns in serum to identify ovarian cancer. *Lancet* **359**, 572–577.
- Ray, S. and Mallick, B. (2006). Functional clustering by Bayesian Wavelet methods. *Journal of the Royal Statistical Society Series B Statistical Methodology* **68**, 305–332.
- Robinson, G. K. (1991). That BLUP is a good thing: The estimation of random effects. *Statistical Science* **6**, 15–32.
- van de Wiel, M. A., Picard, F., van Wieringen, W. N., and Ylstra, B. (2011). Preprocessing and downstream analysis of microarray DNA copy number profiles. *Briefings in Bioinformatics* **12**, 10–21.
- Van Wieringen, W. N., Van De Wiel, M. A., and Ylstra, B. (2008). Weighted clustering of called array CGH data. *Biostatistics* **9**, 484–500.
- Yang, C., He, Z., and Yu, W. (2009). Comparison of public peak detection algorithms for MALDI mass spectrometry data analysis. *BMC Bioinformatics* **10**, 1–13.
- Zhang, J. and Walter, G. (1994). A wavelet-based KL-like expansion for wide sense stationary random processes. *Signal Processing, IEEE Transactions on* **42**, 1737–1745.

Received July 2011. Revised July 2011.

Accepted August 2012.