

# Component-based regularisation of multivariate generalised linear mixed models

Jocelyn Chauvet, Catherine Trottier, Xavier Bry

### ▶ To cite this version:

Jocelyn Chauvet, Catherine Trottier, Xavier Bry. Component-based regularisation of multivariate generalised linear mixed models. Journal of Computational and Graphical Statistics, In press, 10.1080/10618600.2019.1598870. hal-02064508

# HAL Id: hal-02064508 https://hal.science/hal-02064508

Submitted on 9 Aug 2019

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers. L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Component–based regularisation of multivariate generalised linear mixed models

Jocelyn Chauvet \* Catherine Trottier<sup>†\*</sup> and Xavier Bry<sup>\*</sup>

#### Abstract

We address the component-based regularisation of a multivariate Generalised Linear Mixed Model (GLMM) in the framework of grouped data. A set  $\mathbf{Y}$  of random responses is modelled with a multivariate GLMM, based on a set  $\mathbf{X}$  of explanatory variables, a set  $\mathbf{A}$  of additional explanatory variables, and random effects to introduce the within-group dependence of observations. Variables in  $\mathbf{X}$  are assumed many and redundant so that regression demands regularisation. This is not the case for  $\mathbf{A}$ , which contains few and selected variables. Regularisation is performed building an appropriate number of orthogonal components that both contribute to model  $\mathbf{Y}$  and capture relevant structural information in  $\mathbf{X}$ . To estimate the model, we propose to maximise a criterion specific to the Supervised Component-based Generalised Linear Regression (SCGLR) within an adaptation of Schall's algorithm. This extension of SCGLR is tested on both simulated and real grouped data, and compared to ridge and LASSO regularisations. Supplementary material for this article is available online.

*Keywords:* generalised linear regression, supervised components, random effects, structural relevance.

<sup>\*</sup>IMAG, Univ Montpellier, CNRS, Montpellier, France. jocelyn.chauvet@umontpellier.fr; xavier.bry@umontpellier.fr

<sup>&</sup>lt;sup>†</sup>Univ Paul–Valéry Montpellier 3, Montpellier, France. catherine.trottier@univ-montp3.fr

# 1 Introduction

In the framework of regression models on a large number of explanatory variables with redundancies and collinearities, the search for a reduced number of relevant dimensions to model responses has been an ongoing research over the last decades. In particular, the case where the explanatory variables outnumber the observations tends to be a new standard. Generalised Linear Models (GLMs) are the most widely used regression models, because they are easy to interpret and address a very large scope of applications with a variety of response distributions. For instance, Epidemiology, Biology and Social Sciences need to model binary outcomes, count data and survival times. All these fields often have to deal with both grouped data and multivariate responses combining variables of different types (e.g. one binary and another Poisson). In this work, we particularly aim at modelling abundances of several tree genera on plots of land grouped into forest concessions, using multiple redundant explanatory variables.

As far as dimension-reduction is concerned, two main approaches have been developed. The first one is variable-selection, whereas the second one builds components, i.e. linear combinations of the explanatory variables, which synthesise the useful part of their information. As far as variable-selection is concerned, the most popular method is currently the LASSO, introduced by Tibshirani (1996), which combines the likelihood with a penalty based on the  $L_1$ -norm of the coefficient vector. LASSO is one of the penalty-based regularisation methods, as are also ridge (Hoerl and Kennard, 1970) and elastic-net (Zou and Hastie, 2005). This LASSO selection approach has proved efficient to explain the phenomenon of interest when some of the explanatory variables are the "true" ones, surrounded by a high number of irrelevant others. Nevertheless, it may be very unstable and helpless when the true explanatory dimensions are latent and indirectly measured through highly correlated proxies. This is where the component-based approach turns out to be useful. Bry et al. (2013) have developed a new methodology named Supervised Component-based Generalised Linear Regression (SCGLR), later extended and refined in Bry et al. (2014, 2016, 2018). As in any PLS-type method, the construction of components in SCGLR is guided both by the correlation–structure of variables in the explanatory space and by the

prediction quality of the responses. Nevertheless, unlike PLS, SCGLR involves a general and flexible criterion allowing to specify the type of structure components are wanted to align with in the explanatory space (e.g. variable bundles, principal components, other subspaces). Moreover, SCGLR searches for explanatory directions common to multiple responses with probability distributions in the exponential family, each response being entitled to their own distribution. The current SCGLR method is implemented in the R package SCGLR (Cornu et al., 2018) available at https://CRAN.R-project.org/package=SCGLR and https://github.com/SCnext/SCGLR.

In the present work, we aim at modelling responses with a repeated or grouped design. For this purpose, the use of mixed models with random effects is widespread. Research on variance-component estimation in Generalised Linear Mixed Models (GLMMs) has been very active since the 1980s. For the most general distribution assumptions in such models, parameter estimation faces the intractability of the likelihood expressed as an integral with respect to the random effects. Several numerical approximations of the integral have been proposed: Gaussian quadrature (Anderson and Aitkin, 1985) or adaptive versions of it (Pinheiro and Bates, 1995), Laplace approximation leading to the definition of the penalised quasi-likelihood (Breslow and Clayton, 1993) or modified versions of it (Shun and McCullagh, 1995). An alternative to this type of analytic approximation is a stochastic approximation of the integral calculation via MCMC techniques. In this approach, Zeger and Karim (1991) described an approximate Gibbs sampling for GLMMs, which was extended by Clayton (1996) to more general Metropolis–Hastings algorithms. In parallel, McCulloch (1997) developed the Monte Carlo EM algorithm where the expectation is computed numerically through a Monte Carlo approximation, after generating random effects with a Metropolis–Hastings sampler. Mention can also be made of the recent work by Knudson (2016): her strategy is to approximate the entire likelihood function using random effects simulated from a parametrised importance sampling distribution depending on the data. Unfortunately, these different approaches are not necessarily suitable for the same types of random effect designs (one-dimensional random effect, embedded random effects, etc). In the wake of the first type of approximations, we here adopt the "Joint-Maximisation" strategy (McCulloch, 1997), as introduced for instance by Schall (1991). The model is

iteratively linearised conditional on the random effects and variance components are then estimated using adapted linear mixed models methods. This strategy can be used for any random effect design and is less computationally intensive than Monte Carlo methods. Moreover, it provides us with a linear setting more suitable for the computation of components. Once the components calculated, model parameters can be estimated using any of the aforementioned strategies (see Section 7).

Modelling grouped responses through a GLMM with a large number of explanatory variables is the focus of this paper. The need for dimension-reduction and regularisation has to accommodate the presence of random effects in the model, but our main purpose still remains to investigate the explanatory structure and link it to interpretable dimensions. For Gaussian responses, Eliot et al. (2011) proposed to extend the ridge regression to Linear Mixed Models (LMMs). Based on a penalised complete log-likelihood, the adaptation of the Expectation–Maximisation algorithm they suggest includes a new step to find the best shrinkage parameter using a generalised cross-validation scheme at each iteration. More recently, Scheldorfer et al. (2014) — and also Groll and Tutz (2014) — proposed an  $L_1$ penalised algorithm for fitting a high-dimensional GLMM, using Laplace approximation and an efficient coordinate gradient descent. In this work, we combine Schall's iterative model linearisation with regularisation at each step. However, we do not use a penalty on the coefficient vector's norm — as proposed by Zhang et al. (2017) within the framework of multivariate count data. We rather propose to combine dimension-reduction and predictor-regularisation using supervised components aligning on the most predictive and interpretable directions in the explanatory space.

The paper is organised as follows. In Section 2, we formalise the model and set the main notations used throughout the paper. In Section 3, we present the key features of SCGLR. Section 4 designs an extension of this methodology to mixed models, and particularly to grouped data. In Section 5, our extended method "mixed–SCGLR" is evaluated on simulations and compared to ridge– and LASSO–based regularisations. Finally, in order to highlight the power of mixed–SCGLR in terms of model interpretation, Section 6 presents an application to real data in the Poisson case.

### 2 Model definition and notations

In the framework of a multivariate GLMM, we consider q response-vectors  $y_1, \ldots, y_q$ forming matrix  $Y_{n \times q}$ , to be explained by two categories of explanatory variables. The first category consists of few weakly correlated variables  $A_{n \times r} = [a_1 | \ldots | a_r]$ . These variables are assumed to be interesting per se and their marginal effects need to be precisely quantified. The second category consists of abundant and highly correlated variables  $X_{n \times p} = [x_1 | \ldots | x_p]$  considered as proxies to latent dimensions which must be found and interpreted. Since explanatory variables in A are few, non-redundant and of interest, they are kept as such in the model. By contrast, X may contain several unknown structurally relevant dimensions K < p important to model and predict Y, how many we do not know. X is thus to be searched for an appropriate number of orthogonal components that both capture relevant structural information in X and contribute to model Y.

This work addresses grouped data: the *n* observations form *N* groups. Within each group, observations are not assumed independent. For each response  $y_k$ , a *N*-level random effect  $\xi_k$  is used to model the dependence of observations within each group. Hence, each  $y_k$  is modelled with a GLMM assuming a conditional distribution from the exponential family.

#### Notations and conventions

- ▶ All variables (namely the  $a_i$ 's,  $x_j$ 's and  $y_k$ 's) will be identified with *n*-vectors.
- ▶ We will use bold lowercase letters for vectors (e.g.  $\boldsymbol{u}$ ) and bold capital letters for matrices (e.g.  $\boldsymbol{M}$ ).
- ▶ M being any matrix,  $M^{\mathrm{T}}$  denotes the transpose of M.
- ▶  $I_n$  denotes the identity matrix of size n.
- ▶  $\mathbf{1}_m$  denotes the all-ones vector of size m.
- ▶ Let  $\boldsymbol{u}$  and  $\boldsymbol{v}$  be non-zero vectors in  $\mathbb{R}^d$  and let  $\boldsymbol{M}$  be a symmetric positive definite matrix of size  $d \times d$ . Then  $\langle \boldsymbol{u} | \boldsymbol{v} \rangle_{\boldsymbol{M}} = \boldsymbol{u}^{\mathrm{T}} \boldsymbol{M} \boldsymbol{v}$  refers to the Euclidean scalar product of  $\boldsymbol{u}$  and  $\boldsymbol{v}$  with respect to metric  $\boldsymbol{M}$ . The cosine of the angle between  $\boldsymbol{u}$  and  $\boldsymbol{v}$  with respect to  $\boldsymbol{M}$  is given by  $\cos_{\boldsymbol{M}}(\boldsymbol{u}, \boldsymbol{v}) = \frac{\langle \boldsymbol{u} | \boldsymbol{v} \rangle_{\boldsymbol{M}}}{\|\boldsymbol{u}\|_{\boldsymbol{M}} \|\boldsymbol{v}\|_{\boldsymbol{M}}}$ , where  $\|\boldsymbol{u}\|_{\boldsymbol{M}} = \sqrt{\langle \boldsymbol{u} | \boldsymbol{u} \rangle_{\boldsymbol{M}}}$ .

- ▶ The space spanned by vectors  $u_1, \ldots, u_h$  is denoted by span  $\{u_1, \ldots, u_h\}$ . U being any matrix, span  $\{U\}$  refers to the space spanned by the column–vectors of U.
- ▶ Let  $\mathbb{R}^n$  be endowed with metric W and let Z be a matrix of size  $n \times p$ . Then  $\Pi^W_{\text{span}\{Z\}}$  refers to the W-orthogonal projector onto span  $\{Z\}$ . Let b be a vector in  $\mathbb{R}^n$ . The cosine of the angle between b and span  $\{Z\}$  with respect to W is defined by  $\cos_W(b, \text{span}\{Z\}) = \cos_W(b, \Pi^W_{\text{span}\{Z\}}b)$ .

# **3** SCGLR with additional explanatory variables

In this section we consider the situation where each  $y_k$  is modelled with a GLM (without random effect). For the sake of simplicity, we focus on the single-component SCGLR (K = 1). Section 3.1 briefly recalls some standards for univariate GLMs. Section 3.2 defines the linear predictors considered in the SCGLR methodology, in a multivariate GLM framework with additional explanatory variables. Finally, Section 3.3 introduces the criterion SCGLR maximises to compute the component.

#### 3.1 Notations and main features of univariate GLMs

We refer the reader to McCullagh and Nelder (1989) for a thorough overview of GLMs. This section is only intended to recall the classical iterative scheme performing maximum likelihood (ML) estimation. Let X denote the  $n \times p$  matrix of explanatory variables and  $\beta$  the *p*-dimensional parameter vector. At iteration t + 1, the Fisher Scoring Algorithm (FSA) for ML estimation calculates

$$\boldsymbol{\beta}^{[t+1]} = \left(\boldsymbol{X}^{\mathrm{T}} \boldsymbol{W}^{[t]} \boldsymbol{X}\right)^{-1} \boldsymbol{X}^{\mathrm{T}} \boldsymbol{W}^{[t]} \boldsymbol{z}^{[t]}, \qquad (1)$$

where  $\boldsymbol{z}^{[t]}$  and  $\boldsymbol{W}^{[t]}$  respectively denote the classical working variable and the associated weight matrix at iteration t. As pointed out by Nelder and Wedderburn (1972), update (1) may be interpreted as a weighted least squares step in the linearised model  $\mathcal{M}^{[t]}$  defined by

$$\mathcal{M}^{[t]}: \begin{vmatrix} \boldsymbol{z}^{[t]} = \boldsymbol{X}\boldsymbol{\beta} + \boldsymbol{\zeta}^{[t]} \\ \text{with: } \mathbb{E}\left(\boldsymbol{\zeta}^{[t]}\right) = \boldsymbol{0} \text{ and } \mathbb{V}\left(\boldsymbol{\zeta}^{[t]}\right) = \boldsymbol{W}^{[t]^{-1}}. \end{aligned}$$
(2)

#### **3.2** Linear predictors for SCGLR with multiple responses

We are now considering a multivariate GLM (Fahrmeir and Tutz, 1994). In this context, SCGLR searches for a component common to all the  $y_k$ 's. This component will be denoted f and its *p*-dimensional loading-vector will be denoted u, so that f = Xu. The linear predictor associated with response-vector  $y_k$  then writes

$$\boldsymbol{\eta}_{\boldsymbol{k}} = (\boldsymbol{X}\boldsymbol{u})\,\boldsymbol{\gamma}_{\boldsymbol{k}} + \boldsymbol{A}\boldsymbol{\delta}_{\boldsymbol{k}},\tag{3}$$

where  $\gamma_k$  and  $\delta_k$  are the regression parameters associated respectively with component fand additional explanatory variables A. f being common to all the  $y_k$ 's, predictors are collinear in their X-part. For identification purposes, we impose  $u^T M^{-1} u = 1$ , where M may so far be any  $p \times p$  symmetric positive definite matrix. Let us note  $y_{k,i}$  the *i*-th observation of the *k*-th response-vector and  $H = \{\eta_{k,i} \mid 1 \leq k \leq q, 1 \leq i \leq n\}$  the predictor set. We assume that the q responses are independent conditional on f, and that the *n* observations are independent. The log-density then writes

$$\ell\left(\boldsymbol{Y}|\boldsymbol{H}
ight) = \sum_{i=1}^{n}\sum_{k=1}^{q}\ell_{k}\left(y_{k,i}|\eta_{k,i}
ight),$$

where  $\ell_k$  is the log-density of the *k*-th response, conditional on its linear predictor. As a result,  $\boldsymbol{z_k}$  being the working variable associated with  $\boldsymbol{y_k}$  and  $\boldsymbol{W_k}^{-1}$  its variance matrix, the corresponding linearised model derived from the FSA at iteration *t* is

$$\mathcal{M}_{k}^{[t]}: \begin{vmatrix} \boldsymbol{z}_{k}^{[t]} = (\boldsymbol{X}\boldsymbol{u}) \, \gamma_{k} + \boldsymbol{A}\boldsymbol{\delta}_{k} + \boldsymbol{\zeta}_{k}^{[t]} \\ \text{with: } \mathbb{E}\left(\boldsymbol{\zeta}_{k}^{[t]}\right) = 0 \text{ and } \mathbb{V}\left(\boldsymbol{\zeta}_{k}^{[t]}\right) = \boldsymbol{W}_{k}^{[t]^{-1}}. \end{aligned}$$
(4)

Although linearised models (2) and (4) seem very similar, (4) is no longer linear, owing to the product  $\boldsymbol{u}\gamma_k$ . An alternate version of the FSA must therefore be used:

- (i) Given current values of all the  $\gamma_k$ 's and  $\delta_k$ 's, a new loading-vector  $\boldsymbol{u}$  is obtained by solving an SCGLR-specific program (see Section 3.3 for details).
- (ii) Given a current value of  $\boldsymbol{u}$ , each  $\boldsymbol{z}_{\boldsymbol{k}}$  is regressed independently on  $[\boldsymbol{X}\boldsymbol{u} \mid \boldsymbol{A}]$  with respect to weight matrix  $\boldsymbol{W}_{\boldsymbol{k}}$ , yielding new regression parameters  $\gamma_k$  and  $\boldsymbol{\delta}_{\boldsymbol{k}}$ .

# 3.3 Calculating the component maximising an SCGLR–specific criterion

For an easier reading of this part, we omit the [t] index. For each  $k \in \{1, \ldots, q\}$ , consider model  $\mathcal{M}_k$  endowed with weight matrix  $W_k$ . As suggested in Bry et al. (2013), the best loading-vector in the weighted least-squares sense would be the solution of

$$\min_{\boldsymbol{u}:\,\boldsymbol{u}^{\mathrm{T}}\boldsymbol{M}^{-1}\boldsymbol{u}=1} \sum_{k=1}^{q} \left\|\boldsymbol{z}_{\boldsymbol{k}} - \Pi_{\mathrm{span}\{\boldsymbol{X}\boldsymbol{u},\boldsymbol{A}\}}^{\boldsymbol{W}_{\boldsymbol{k}}}\boldsymbol{z}_{\boldsymbol{k}}\right\|_{\boldsymbol{W}_{\boldsymbol{k}}}^{2} \iff \max_{\boldsymbol{u}:\,\boldsymbol{u}^{\mathrm{T}}\boldsymbol{M}^{-1}\boldsymbol{u}=1} \sum_{k=1}^{q} \left\|\Pi_{\mathrm{span}\{\boldsymbol{X}\boldsymbol{u},\boldsymbol{A}\}}^{\boldsymbol{W}_{\boldsymbol{k}}}\boldsymbol{z}_{\boldsymbol{k}}\right\|_{\boldsymbol{W}_{\boldsymbol{k}}}^{2}.$$

The maximisation program also writes  $\max_{\boldsymbol{u}: \boldsymbol{u}^{\mathrm{T}} \boldsymbol{M}^{-1} \boldsymbol{u}=1} \psi_{\boldsymbol{A}}(\boldsymbol{u})$ , where

$$\psi_{\boldsymbol{A}}(\boldsymbol{u}) = \sum_{k=1}^{q} \left\| \boldsymbol{z}_{\boldsymbol{k}} \right\|_{\boldsymbol{W}_{\boldsymbol{k}}}^{2} \cos^{2}_{\boldsymbol{W}_{\boldsymbol{k}}} \left( \boldsymbol{z}_{\boldsymbol{k}}, \operatorname{span} \left\{ \boldsymbol{X} \boldsymbol{u}, \boldsymbol{A} \right\} \right)$$
$$= \sum_{k=1}^{q} \left\| \boldsymbol{z}_{\boldsymbol{k}} \right\|_{\boldsymbol{W}_{\boldsymbol{k}}}^{2} \cos^{2}_{\boldsymbol{W}_{\boldsymbol{k}}} \left( \boldsymbol{z}_{\boldsymbol{k}}, \Pi_{\operatorname{span}\left\{ \boldsymbol{X} \boldsymbol{u}, \boldsymbol{A} \right\}}^{\boldsymbol{W}_{\boldsymbol{k}}} \boldsymbol{z}_{\boldsymbol{k}} \right).$$
(5)

Now,  $\psi_{A}$  is a mere goodness-of-fit (GoF) measure that does not take into account the closeness of component  $\boldsymbol{f} = \boldsymbol{X}\boldsymbol{u}$  to interpretable directions in  $\boldsymbol{X}$ . The GoF measure,  $\psi_{A}$ , must therefore be combined with a measure  $\phi$  of structural relevance (SR).

Assume matrix  $\boldsymbol{X}$  consists of p standardised numeric variables. Consider a weight system  $\boldsymbol{\omega} = \{\omega_1, \ldots, \omega_p\}$  — e.g.  $\omega_j = \frac{1}{p} \forall j \in \{1, \ldots, p\}$  — reflecting the a priori relative importance of variables. Also consider a weight matrix  $\boldsymbol{P}$  — e.g.  $\boldsymbol{P} = \frac{1}{n}\boldsymbol{I}_n$  — reflecting the a priori relative importance of observations. We define the most structurally relevant loading-vector as the solution of

$$\max_{\boldsymbol{u}:\,\boldsymbol{u}^{\mathrm{T}}\boldsymbol{M}^{-1}\boldsymbol{u}=1} \phi(\boldsymbol{u}),$$

where

$$\phi\left(\boldsymbol{u}\right) = \left[\sum_{j=1}^{p} \omega_{j} \left(\left\langle \boldsymbol{X}\boldsymbol{u} \,|\, \boldsymbol{x}_{\boldsymbol{j}} \right\rangle_{\boldsymbol{P}}^{2}\right)^{l}\right]^{\frac{1}{l}} = \left[\sum_{j=1}^{p} \omega_{j} \left(\boldsymbol{u}^{\mathrm{T}} \boldsymbol{X}^{\mathrm{T}} \boldsymbol{P} \boldsymbol{x}_{\boldsymbol{j}} \boldsymbol{x}_{\boldsymbol{j}}^{\mathrm{T}} \boldsymbol{P} \boldsymbol{X} \,\boldsymbol{u}\right)^{l}\right]^{\frac{1}{l}}, \ l \ge 1, \quad (6)$$

for the scalar product is commutative. Formula (6) is in fact a particular case of the SR criterion proposed by Bry and Verron (2015); Bry et al. (2016). It can be viewed as a generalised average version of the usual dual PCA criterion:  $\sum_{j=1}^{p} \cos^{2}_{P} (Xu, x_{j}) =$ 

 $\sum_{j=1}^{p} \langle X \boldsymbol{u} | \boldsymbol{x}_{j} \rangle_{\boldsymbol{P}}^{2}$ . For  $\boldsymbol{M} = (\boldsymbol{X}^{T} \boldsymbol{P} \boldsymbol{X})^{-1}$ , (6) is called "Variable–Powered Inertia" (VPI). It should be stressed that for  $\boldsymbol{X}^{T} \boldsymbol{P} \boldsymbol{X}$  to be invertible,  $\boldsymbol{X}$  must be a column full rank matrix. In case of strict collinearities within  $\boldsymbol{X}$ , as it always happens in high–dimensional settings, we replace  $\boldsymbol{X}$  with the matrix  $\boldsymbol{C}$  of its principal components associated with non–zero eigenvalues. The component is then sought as  $\boldsymbol{f} = \boldsymbol{C}\boldsymbol{u}$ . We have  $\boldsymbol{C} = \boldsymbol{X}\boldsymbol{V}$ , where  $\boldsymbol{V}$  is the matrix of corresponding unit-eigenvectors. Then,  $\boldsymbol{f} = \boldsymbol{C}\boldsymbol{u} = \boldsymbol{X}\tilde{\boldsymbol{u}}$  with  $\tilde{\boldsymbol{u}} = \boldsymbol{V}\boldsymbol{u}$ . Bry et al. (2018) show that among all loading–vectors  $\boldsymbol{t}$  such that  $\boldsymbol{X}\boldsymbol{t} = \boldsymbol{f}, \, \tilde{\boldsymbol{u}}$  is that which has the minimum  $L_2$ –norm.

Tuning parameter l allows to draw component towards more (greater l) or less (smaller l) local bundles of correlated variables, as depicted on Figure 1 in the particular instance of four coplanar variables. Informally, a bundle is a set of variables correlated "enough" to be viewed as proxies to the same latent dimension. The notion of bundle is flexible, and parameter l tunes the level of within-bundle correlation to be considered: the higher the correlation, the more local the bundle. Overall, taking l = 1 draws the components towards global structural directions (namely the principal components) while taking l higher leads to more local ones (ultimately, the variables themselves). The goal is to focus on the most interpretable directions.

Finally, let  $s \in [0, 1]$  be a parameter tuning the importance of the SR relative to the GoF. SCGLR attempts a trade-off between (5) and (6) by solving

$$\max_{\boldsymbol{u}:\,\boldsymbol{u}^{\mathrm{T}}\boldsymbol{M}^{-1}\boldsymbol{u}=1} \left[\phi\left(\boldsymbol{u}\right)\right]^{s} \left[\psi_{\boldsymbol{A}}\left(\boldsymbol{u}\right)\right]^{1-s}$$

or equivalently

$$\max_{\boldsymbol{\mu}: \boldsymbol{u}^{\mathrm{T}} \boldsymbol{M}^{-1} \boldsymbol{u} = 1} s \log \left[ \phi \left( \boldsymbol{u} \right) \right] + (1 - s) \log \left[ \psi_{\boldsymbol{A}} \left( \boldsymbol{u} \right) \right].$$
(7)

More detail can be found in Bry et al. (2018).

### 4 Extension to mixed models

We now propose to extend SCGLR to mixed models. This extension will be called "mixed– SCGLR". A particular focus is placed on grouped data, for which the independence as-



Figure 1: Polar representation of the VPI according to the value of l in the elementary case of four coplanar variables,  $\mathbf{x_1}, \mathbf{x_2}, \mathbf{x_3}, \mathbf{x_4}$ , with  $\omega_j = \frac{1}{4} \forall j \in \{1, 2, 3, 4\}$ . Loading-vector  $\mathbf{u}$ is identified with complex number  $e^{i\theta}$ , where  $\theta \in [0, 2\pi)$ . Curves  $z_l(\theta) := [\phi(e^{i\theta})]^l e^{i\theta}$  are graphed for  $l \in \{1, 2, 4, 10, 50\}$ . The intersection of curve  $z_l$  with  $\mathbf{f} = \mathbf{X}\mathbf{u}$  has a radius equal to  $[\phi(e^{i\theta})]^l$ . The red line is the direction of maximum for l = 1, which is in fact the first principal component. These four variables are then regarded as a unique bundle. By contrast, the blue lines represent the two directions of maximum for l = 4. The variables are then seen as two bundles containing two variables each. Finally, when l = 50, each variable is considered a bundle in itself.

sumption of observations is no longer valid. The within-group dependence of each response is modelled with a random group-effect. Consequently, each  $y_k$  is modelled with a GLMM. As in SCGLR, the responses are assumed to be independent conditional on the components. Section 4.1 presents the single-component mixed-SCGLR method. The underlying algorithm is given in Section 4.2. Considering only one component is generally not enough to explain the responses making it necessary to search for K explanatory components, with  $1 \leq K \leq \operatorname{rank}(\mathbf{X})$ . The way in which we extract higher rank components is explained in Section 4.3.

#### 4.1 First component

The random group–effect is assumed different across responses. This leads to q random–effect vectors  $\boldsymbol{\xi}_1, \ldots, \boldsymbol{\xi}_q$ , which are assumed independent and normally distributed:

$$\forall k \in \{1, \ldots, q\}, \quad \boldsymbol{\xi}_{\boldsymbol{k}} \stackrel{\text{ind.}}{\sim} \mathcal{N}_{N}(\boldsymbol{0}, \boldsymbol{D}_{\boldsymbol{k}}),$$

where N denotes the number of groups. In this paper, variance components models will be considered. We assume  $D_k = \sigma_k^2 I_N$ , where  $\sigma_k^2$  is the group variance component associated with response  $y_k$ . Linear predictors involved in mixed–SCGLR are expressed as

$$\forall k \in \{1, \dots, q\}, \quad \boldsymbol{\eta}_{\boldsymbol{k}}^{\boldsymbol{\xi}} = (\boldsymbol{X}\boldsymbol{u})\gamma_k + \boldsymbol{A}\boldsymbol{\delta}_{\boldsymbol{k}} + \boldsymbol{U}\boldsymbol{\xi}_{\boldsymbol{k}}, \tag{8}$$

where U is the known random effects' design matrix. Predictor  $\eta_k^{\xi}$  epitomises the way we capture the dependence between outcomes. Indeed, as component f = Xu does not depend on k, it captures a structural dependence between the various  $y_k$ 's. By contrast, the random effect  $\xi_k$  models the within-group stochastic dependence of outcomes forming response-vector  $y_k$ .

Recall that the distribution of the data conditional on the random effects is supposed to belong to the exponential family. The FSA was adapted by Schall (1991) to the GLMM dependence structure. The key idea is to extend Schall's algorithm to the component-based predictors in (8).

#### 4.1.1 Linearisation step

Let  $g_k$  denote the link function for response  $\boldsymbol{y}_k$ ,  $g'_k$  its first derivative and  $\boldsymbol{\mu}_k^{\boldsymbol{\xi}}$  the conditional expectation (i.e.  $\boldsymbol{\mu}_k^{\boldsymbol{\xi}} := \mathbb{E}(\boldsymbol{y}_k | \boldsymbol{\xi}_k)$ ). The working variable associated with  $y_{k,i}$  is calculated through

$$z_{k,i}^{\xi} = g_k \left( \mu_{k,i}^{\xi} \right) + \left( y_{k,i} - \mu_{k,i}^{\xi} \right) g'_k \left( \mu_{k,i}^{\xi} \right) \\ = \eta_{k,i}^{\xi} + e_{k,i}, \quad \text{where} \quad e_{k,i} = \left( y_{k,i} - \mu_{k,i}^{\xi} \right) g'_k \left( \mu_{k,i}^{\xi} \right).$$

In view of the conditional independence assumption, the conditional variance matrix for  $z_k^{\xi}$  is

$$\operatorname{Var}\left(\boldsymbol{z}_{\boldsymbol{k}}^{\boldsymbol{\xi}} \,|\, \boldsymbol{\xi}_{\boldsymbol{k}}\right) = \boldsymbol{W}_{\boldsymbol{k}}^{\boldsymbol{\xi}^{-1}} = \operatorname{Diag}\left(\left[g_{\boldsymbol{k}}'\left(\boldsymbol{\mu}_{\boldsymbol{k},i}^{\boldsymbol{\xi}}\right)\right]^{2} a_{\boldsymbol{k},i}(\phi_{\boldsymbol{k}}) \,v_{\boldsymbol{k}}\left(\boldsymbol{\mu}_{\boldsymbol{k},i}^{\boldsymbol{\xi}}\right)\right)_{i=1,\dots,n},$$

where  $a_{k,i}$  and  $v_k$  are known functions, and  $\phi_k$  is the dispersion parameter related to  $y_k$ . At iteration t, the conditional linearised model for working vector  $\boldsymbol{z}_k^{\boldsymbol{\xi}}$  is then defined by

$$\mathcal{M}_{k}^{\boldsymbol{\xi}^{[t]}}: \begin{vmatrix} \boldsymbol{z_{k}^{\boldsymbol{\xi}^{[t]}}} = (\boldsymbol{X}\boldsymbol{u}) \, \gamma_{k} + \boldsymbol{A}\boldsymbol{\delta}_{\boldsymbol{k}} + \boldsymbol{U}\boldsymbol{\xi}_{\boldsymbol{k}} + \boldsymbol{e}_{\boldsymbol{k}}^{[t]} \\ \text{with: } \mathbb{E}\left(\boldsymbol{e}_{\boldsymbol{k}}^{[t]} \, | \, \boldsymbol{\xi}_{\boldsymbol{k}}\right) = 0 \text{ and } \mathbb{V}\left(\boldsymbol{e}_{\boldsymbol{k}}^{[t]} \, | \, \boldsymbol{\xi}_{\boldsymbol{k}}\right) = \boldsymbol{W}_{\boldsymbol{k}}^{\boldsymbol{\xi}^{-1}[t]}. \end{aligned}$$
(9)

Besides the variance component estimation, an alternated estimation step has to be developed (as aforementioned in Section 3.2) to deal with the non–linearity of (9).

#### 4.1.2 Estimation step

Calculating the component: Given current values of all the  $\gamma_k$ 's,  $\delta_k$ 's,  $\xi_k$ 's and  $\sigma_k^2$ 's, a new component f = Xu is calculated by solving a (7)-type program. However, (5) has to be adapted to conditional linearised models  $\mathcal{M}_k^{\xi}$ 's, involving weight matrices  $W_k^{\xi}$ 's. The appropriate goodness-of-fit measure is

$$\psi_{\boldsymbol{A}}(\boldsymbol{u}) = \sum_{k=1}^{q} \left\| \boldsymbol{z}_{\boldsymbol{k}}^{\boldsymbol{\xi}} \right\|_{\boldsymbol{W}_{\boldsymbol{k}}^{\boldsymbol{\xi}}}^{2} \cos^{2}_{\boldsymbol{W}_{\boldsymbol{k}}^{\boldsymbol{\xi}}} \left( \boldsymbol{z}_{\boldsymbol{k}}^{\boldsymbol{\xi}}, \operatorname{span}\left\{ \boldsymbol{X}\boldsymbol{u}, \boldsymbol{A} \right\} \right).$$
(10)

Estimating the regression parameters and variance–components: Given a current value of component f, we apply Schall's method with the linear predictors given in (8). New values of parameters  $\gamma_k$  and  $\delta_k$  as well as new prediction  $\xi_k$  are obtained by solving the following Henderson system (Henderson, 1975):

$$egin{pmatrix} f^{^{\mathrm{T}}}W_k^{\xi}f & f^{^{\mathrm{T}}}W_k^{\xi}A & f^{^{\mathrm{T}}}W_k^{\xi}U \ A^{^{\mathrm{T}}}W_k^{\xi}f & A^{^{\mathrm{T}}}W_k^{\xi}A & A^{^{\mathrm{T}}}W_k^{\xi}U \ U^{^{\mathrm{T}}}W_k^{\xi}f & U^{^{\mathrm{T}}}W_k^{\xi}A & U^{^{\mathrm{T}}}W_k^{\xi}U+D_k^{-1} \end{pmatrix} egin{pmatrix} \gamma_k \ \delta_k \ \xi_k \end{pmatrix} = egin{pmatrix} f^{^{\mathrm{T}}}W_k^{\xi}z_k^{\xi} \ U^{^{\mathrm{T}}}W_k^{\xi}z_k^{\xi} \ U^{^{\mathrm{T}}}W_k^{\xi}z_k^{\xi} \end{pmatrix} \end{split}$$

Finally, as mentioned by Schall (1991), given prediction  $\hat{\xi}_k$  for  $\xi_k$ , the update of the ML estimation of variance component  $\sigma_k^2$  is

$$\sigma_k^2 \longleftarrow \frac{\widehat{\boldsymbol{\xi}_k}^{\mathrm{T}} \widehat{\boldsymbol{\xi}_k}}{N - \frac{1}{\sigma_k^2} \operatorname{Trace}\left[ \left( \boldsymbol{U}^{\mathrm{T}} \boldsymbol{W}_k^{\boldsymbol{\xi}} \boldsymbol{U} + \boldsymbol{D}_k^{-1} \right)^{-1} \right]}.$$

#### 4.2 The algorithm

The conditional linearised models considered at iteration t are given by (9). Algorithm 1 describes the (t + 1)-th iteration of the single-component mixed-SCGLR. It is repeated until stability of parameters is reached.

Step 1: Computing the component. Set  $\boldsymbol{u}^{[t+1]} = \underset{\boldsymbol{u}: \boldsymbol{u}^{\mathrm{T}} \boldsymbol{M}^{-1} \boldsymbol{u}=1}{\operatorname{arg max}} \left[ \phi\left(\boldsymbol{u}\right) \right]^{s} \left[ \psi_{\boldsymbol{A}}^{[t]}\left(\boldsymbol{u}\right) \right]^{1-s}, \text{ where } \psi_{\boldsymbol{A}}\left(\boldsymbol{u}\right) \text{ is given by (10) and } \phi\left(\boldsymbol{u}\right) \text{ by (6)}$   $\boldsymbol{f}^{[t+1]} = \boldsymbol{X} \boldsymbol{u}^{[t+1]}$ 

**Step 2: Henderson systems.** For each  $k \in \{1, \ldots, q\}$ , solve the system

$$\begin{pmatrix} \boldsymbol{f}^{[t+1]^{\mathrm{T}}} \boldsymbol{W}_{\boldsymbol{k}}^{\boldsymbol{\xi}[t]} \boldsymbol{f}^{[t+1]} & \boldsymbol{f}^{[t+1]^{\mathrm{T}}} \boldsymbol{W}_{\boldsymbol{k}}^{\boldsymbol{\xi}[t]} \boldsymbol{A} & \boldsymbol{f}^{[t+1]^{\mathrm{T}}} \boldsymbol{W}_{\boldsymbol{k}}^{\boldsymbol{\xi}[t]} \boldsymbol{U} \\ \boldsymbol{A}^{\mathrm{T}} \boldsymbol{W}_{\boldsymbol{k}}^{\boldsymbol{\xi}[t]} \boldsymbol{f}^{[t]} & \boldsymbol{A}^{\mathrm{T}} \boldsymbol{W}_{\boldsymbol{k}}^{\boldsymbol{\xi}[t]} \boldsymbol{A} & \boldsymbol{A}^{\mathrm{T}} \boldsymbol{W}_{\boldsymbol{k}}^{\boldsymbol{\xi}[t]} \boldsymbol{U} \\ \boldsymbol{U}^{\mathrm{T}} \boldsymbol{W}_{\boldsymbol{k}}^{\boldsymbol{\xi}[t]} \boldsymbol{f}^{[t]} & \boldsymbol{U}^{\mathrm{T}} \boldsymbol{W}_{\boldsymbol{k}}^{\boldsymbol{\xi}[t]} \boldsymbol{A} & \boldsymbol{U}^{\mathrm{T}} \boldsymbol{W}_{\boldsymbol{k}}^{\boldsymbol{\xi}[t]} \boldsymbol{U} + \boldsymbol{D}_{\boldsymbol{k}}^{[t]^{-1}} \end{pmatrix} \begin{pmatrix} \boldsymbol{\gamma}_{\boldsymbol{k}} \\ \boldsymbol{\delta}_{\boldsymbol{k}} \\ \boldsymbol{\xi}_{\boldsymbol{k}} \end{pmatrix} = \begin{pmatrix} \boldsymbol{f}^{[t+1]^{\mathrm{T}}} \boldsymbol{W}_{\boldsymbol{k}}^{\boldsymbol{\xi}[t]} \boldsymbol{z}_{\boldsymbol{k}}^{\boldsymbol{\xi}[t]} \\ \boldsymbol{U}^{\mathrm{T}} \boldsymbol{W}_{\boldsymbol{k}}^{\boldsymbol{\xi}[t]} \boldsymbol{z}_{\boldsymbol{k}}^{\boldsymbol{\xi}[t]} \\ \boldsymbol{U}^{\mathrm{T}} \boldsymbol{W}_{\boldsymbol{k}}^{\boldsymbol{\xi}[t]} \boldsymbol{z}_{\boldsymbol{k}}^{\boldsymbol{\xi}[t]} \end{pmatrix} \\ \text{Call } \boldsymbol{\gamma}_{\boldsymbol{k}}^{[t+1]}, \boldsymbol{\delta}_{\boldsymbol{k}}^{[t+1]} \text{ and } \boldsymbol{\xi}_{\boldsymbol{k}}^{[t+1]} \text{ the solutions.} \end{cases}$$

Step 3: Updating variance–component estimates. For each  $k \in \{1, \ldots, q\}$ , compute

$$\sigma_{k}^{2^{[t+1]}} = \frac{\boldsymbol{\xi}_{k}^{[t+1]^{\mathrm{T}}} \, \boldsymbol{\xi}_{k}^{[t+1]}}{N - \frac{1}{\sigma_{k}^{2^{[t]}}} \operatorname{Trace}\left[ \left( \boldsymbol{U}^{\mathrm{T}} \, \boldsymbol{W}_{k}^{\boldsymbol{\xi}^{[t]}} \, \boldsymbol{U} + \boldsymbol{D}_{k}^{[t]^{-1}} \right)^{-1} \right]} \quad \text{and} \quad \boldsymbol{D}_{k}^{[t+1]} = \sigma_{k}^{2^{[t+1]}} \boldsymbol{I}_{N}$$

Step 4: Updating working variables and weighting matrices.

For each  $k \in \{1, \ldots, q\}$ , compute

$$\begin{split} \boldsymbol{\eta}_{\boldsymbol{k}}^{\boldsymbol{\xi}^{[t+1]}} &= \boldsymbol{f}^{[t+1]} \boldsymbol{\gamma}_{k}^{[t+1]} + \boldsymbol{A} \boldsymbol{\delta}_{\boldsymbol{k}}^{[t+1]} + \boldsymbol{U} \boldsymbol{\xi}_{\boldsymbol{k}}^{[t+1]} \\ \boldsymbol{\mu}_{k,i}^{\boldsymbol{\xi}} \stackrel{[t+1]}{=} g_{k}^{-1} \left( \boldsymbol{\eta}_{k,i}^{\boldsymbol{\xi}} \stackrel{[t+1]}{)} \right), \ i = 1, \dots, n \\ \boldsymbol{z}_{k,i}^{\boldsymbol{\xi}} \stackrel{[t+1]}{=} \eta_{k,i}^{\boldsymbol{\xi}} \stackrel{[t+1]}{=} + \left( \boldsymbol{y}_{i}^{\boldsymbol{\xi}} - \boldsymbol{\mu}_{k,i}^{\boldsymbol{\xi}} \stackrel{[t+1]}{)} \right) \boldsymbol{g}_{k}^{\prime} \left( \boldsymbol{\mu}_{k,i}^{\boldsymbol{\xi}} \stackrel{[t+1]}{)} \right), \ i = 1, \dots, n \\ \boldsymbol{W}_{\boldsymbol{k}}^{\boldsymbol{\xi}} \stackrel{[t+1]}{=} \mathbf{Diag} \left( \left\{ \left[ \boldsymbol{g}_{k}^{\prime} \left( \boldsymbol{\mu}_{k,i}^{\boldsymbol{\xi}} \right) \right]^{2} \boldsymbol{a}_{k,i}(\boldsymbol{\phi}_{k}) \, \boldsymbol{v}_{k} \left( \boldsymbol{\mu}_{k,i}^{\boldsymbol{\xi}} \stackrel{[t+1]}{)} \right) \right\}^{-1} \right)_{i=1,\dots,n} \end{split}$$

Incrementing:  $t \leftarrow t+1$ 

Algorithm 1: Iteration of the single-component mixed-SCGLR

#### 4.3 Extracting higher rank components

Let  $F_h = [f_1 | \dots | f_h]$  be the matrix of the first h components, where h < K. An extra component  $f_{h+1}$  must best complement the existing ones plus A, i.e.  $A_h = [F_h | A]$ . So  $f_{h+1}$  must be calculated using  $A_h$  as additional explanatory variables. Moreover, we must impose that  $f_{h+1}$  be orthogonal to  $F_h$ , i.e.  $F_h^T P f_{h+1} = 0$ . Component  $f_{h+1} = X u_{h+1}$  is thus obtained by solving

$$\begin{cases} \max \quad s \log \left[ \phi \left( \boldsymbol{u} \right) \right] + (1 - s) \log \left[ \psi_{\boldsymbol{A}_{\boldsymbol{h}}} \left( \boldsymbol{u} \right) \right] \\ \text{subject to:} \quad \boldsymbol{u}^{\mathrm{T}} \boldsymbol{M}^{-1} \boldsymbol{u} = 1 \text{ and } \boldsymbol{D}_{\boldsymbol{h}}^{\mathrm{T}} \boldsymbol{u} = \boldsymbol{0}, \end{cases}$$
(11)  
where  $\psi_{\boldsymbol{A}_{\boldsymbol{h}}} \left( \boldsymbol{u} \right) = \sum_{k=1}^{q} \left\| \boldsymbol{z}_{\boldsymbol{k}}^{\boldsymbol{\xi}} \right\|_{\boldsymbol{W}_{\boldsymbol{k}}^{\boldsymbol{\xi}}}^{2} \cos^{2}_{\boldsymbol{W}_{\boldsymbol{k}}^{\boldsymbol{\xi}}} \left( \boldsymbol{z}_{\boldsymbol{k}}^{\boldsymbol{\xi}}, \operatorname{span} \left\{ \boldsymbol{X} \boldsymbol{u}, \boldsymbol{A}_{\boldsymbol{h}} \right\} \right) \text{ and } \boldsymbol{D}_{\boldsymbol{h}} = \boldsymbol{X}^{\mathrm{T}} \boldsymbol{P} \boldsymbol{F}_{\boldsymbol{h}}.$ 

In the online Supplementary Material, we give a simple tool to maximise, at least locally, any criterion on the unit sphere: the Projected Iterated Normed Gradient (PING) algorithm. In particular, PING solves (11)-type programs, which give all components of rank h > 1. The rank-one component is computed using the same program with  $A_0 = A$ and  $D_0 = 0$ .

### 5 Comparative results on simulated data

Five simulation studies have been implemented to assess our method. The first one (discussed in Sections 5.1 - 5.4) focuses on LMMs. It compares the performances of mixed–SCGLR, LMM–ridge (Eliot et al., 2011) and GLMM–LASSO (Groll and Tutz, 2014; Schelldorfer et al., 2014). The second simulation (Section 5.5) extends the first one to binary and Poisson outcomes. All simulation studies have been performed using R (R Core Team, 2017). To compute LASSO regressions, we have used the R package glmmLasso (Groll, 2017). The extension of SCGLR to mixed models is available at https://github.com/SCnext/mixedSCGLR. Three additional simulations are presented in the online Supplementary Material. The first one reproduces the simulation scheme of Section 5.5 with binomial and Poisson outcomes. The second one assesses the performance of mixed–SCGLR on a different bundle structure and presents results concerning variance component estimates. The third one deals with high dimensional data.

#### 5.1 Data generation

To generate grouped data, we consider N = 10 groups and R = 10 observations per group (i.e. a total of n = 100 observations). The random effects' design matrix is then  $U = I_N \otimes \mathbf{1}_R$ . Explanatory variables X consist of three independent bundles:  $X_0$  (15 variables),  $X_1$  (10 variables) and  $X_2$  (5 variables). Each explanatory variable is normally simulated with mean 0 and variance 1. Parameter  $\tau \in \{0.1, 0.3, 0.5, 0.7, 0.9\}$  tunes the level of redundancy within each bundle: the correlation matrix of bundle  $X_j$  is

$$\operatorname{cor}\left(\boldsymbol{X}_{\boldsymbol{j}}\right) = \tau \mathbf{1}_{p_{\boldsymbol{j}}} \mathbf{1}_{p_{\boldsymbol{j}}}^{\mathrm{T}} + (1-\tau) \boldsymbol{I}_{p_{\boldsymbol{j}}},$$

where  $p_j$  is the number of variables in  $X_j$ . In order to enable comparison with LASSO and ridge and to focus on regularisation, our simulations do not involve additional explanatory variables (A = 0). Two random responses  $Y = [y_1 | y_2]$  are generated as

$$\begin{cases} \boldsymbol{y}_1 = \boldsymbol{X}\boldsymbol{\beta}_1 + \boldsymbol{U}\boldsymbol{\xi}_1 + \boldsymbol{\varepsilon}_1 \\ \boldsymbol{y}_2 = \boldsymbol{X}\boldsymbol{\beta}_2 + \boldsymbol{U}\boldsymbol{\xi}_2 + \boldsymbol{\varepsilon}_2, \end{cases}$$
(12)

such that  $y_1$  is predicted only by bundle  $X_1$ ,  $y_2$  only by bundle  $X_2$ , and bundle  $X_0$  plays no explanatory role. Our choice for the fixed-effect parameters is

$$\beta_{1} = (\underbrace{0, \dots, 0, 0, 0.3, \dots, 0.3}_{15 \text{ times}}, \underbrace{0.4, \dots, 0.4}_{4 \text{ times}}, \underbrace{0.5, \dots, 0.5}_{3 \text{ times}}, \underbrace{0, \dots, 0}_{5 \text{ times}})^{\mathrm{T}}, \\ \beta_{2} = (\underbrace{0, \dots, 0, 0}_{25 \text{ times}}, \underbrace{0, \dots, 0, \dots, 0}_{$$

Finally, for each  $k \in \{1, 2\}$ , random effect and noise vectors are simulated respectively from

$$\boldsymbol{\xi_k} \sim \mathcal{N}_N\left( \mathbf{0}, \ \sigma_k^2 \, \boldsymbol{I}_N 
ight) \ ext{and} \ \boldsymbol{\varepsilon_k} \sim \mathcal{N}_n\left( \mathbf{0}, \ \omega_k^2 \, \boldsymbol{I}_n 
ight),$$

where  $\sigma_k^2 = \omega_k^2 = 1$ . For each value of  $\tau$ , B = 100 samples are generated according to Model (12).

#### 5.2 Parameter calibration

In order to compare mixed–SCGLR with the ridge and LASSO regressions, we recall the regularisation parameters required by each method. For both LMM–ridge and GLMM–LASSO methods, a unique shrinkage parameter has to be calibrated:  $\lambda_{ridge}$  and  $\lambda_{LASSO}$  respectively. For mixed–SCGLR, three tuning parameters need to be calibrated: the number

of components K and the trade-off parameter s, which are both regularisation parameters, and the bundle-locality parameter l. For greater clarity, the simulation focuses on the behaviour of K and s. As recommended by Bry et al. (2013), we set l = 4. In case-studies, l has to be tuned to maximise the interpretability of components.

For both mixed–SCGLR and GLMM–LASSO, optimal regularisation parameters are obtained through a 5–fold cross–validation, withdrawing 2 observations from each group every time. This could be termed "leave–two–observations–out per group." The data are thus divided into five parts  $\mathcal{P}_1, \ldots, \mathcal{P}_5$ , each  $\mathcal{P}_j$  containing 20 observations, 2 for each of the 10 groups. Let  $y_{k,i}^{(b)}$  be the *i*–th observation of the *k*–th response vector in the *b*–th sample. Let also  $\widehat{y_{k,i(-j)}^{(b)}}$  be the fit for  $y_{k,i}^{(b)}$  with part  $\mathcal{P}_j$  removed. The cross–validation error in the *b*–th sample,  $E^{(b)}$ , is defined as

$$E^{(b)} = \frac{1}{2} \sum_{k=1}^{2} E_k^{(b)}, \tag{13}$$

where

$$E_k^{(b)} = \frac{1}{5} \sum_{j=1}^5 \sqrt{\frac{1}{20} \sum_{i \in \mathcal{P}_j} \left( y_{k,i}^{(b)} - \widehat{y_{k,i(-j)}^{(b)}} \right)^2}.$$

In the *b*-th sample, the optimal number of components  $K^{\star(b)}$ , the trade-off parameter  $s^{\star(b)}$ , and the shrinkage parameter  $\lambda_{\text{LASSO}}^{\star(b)}$  are selected to minimise the cross-validation error (13). We then define

$$s^{\star} = \frac{1}{B} \sum_{b=1}^{B} s^{\star(b)}, \ K^{\star} = \text{mode}\left(\left\{K^{\star(1)}, \dots, K^{\star(B)}\right\}\right) \text{ and } \lambda_{\text{LASSO}}^{\star} = \frac{1}{B} \sum_{b=1}^{B} \lambda_{\text{LASSO}}^{\star(b)}.$$

By contrast, Eliot et al. (2011) suggest to calibrate the ridge parameter at each step of their EM implementation, using the generalised cross-validation. We thus define

$$\lambda_{\rm ridge}^{\star} = \frac{1}{B} \sum_{b=1}^{B} \lambda_{\rm ridge}^{\star(b)}$$

where  $\lambda_{\text{ridge}}^{\star(b)}$  denotes the average of the ridge parameter values obtained over all the iterations of the EM algorithm in the *b*-th sample.

Table 1 summarises the optimal regularisation parameters selected through cross–validation. In both ridge and LASSO, the shrinkage parameter value increases with the

level of redundancy  $\tau$ . Whereas for mixed-SCGLR, when  $\tau$  increases,  $K^*$  decreases towards the true number of predictive variable-bundles: the greater the value of  $\tau$ , the better mixed-SCGLR focuses on the structures in X that contribute to model Y. Moreover, when  $\tau$  increases, the trade-off parameter  $s^*$  increases, meaning that regularisation requires a greater importance of the structural relevance relative to the goodness-of-fit.

	GLMM-LASSO	$\mathbf{LMM}$ -ridge	$\mathbf{mixed}\mathbf{-}\mathbf{SCGLR}$		
	$\operatorname{shrinkage}$	shrinkage	number of	trade-off	
	parameter $\lambda^{\star}_{\text{LASSO}}$	parameter $\lambda_{\text{ridge}}^{\star}$	components $K^{\star}$	parameter $s^{\star}$	
$\tau = 0.1$	65	24	15	0.50	
$\tau = 0.3$	92	54	5	0.58	
$\tau = 0.5$	124	73	3	0.70	
$\tau = 0.7$	163	78	3	0.73	
$\tau = 0.9$	175	85	2	0.80	

Table 1: Optimal regularisation parameter values obtained through cross-validation over 100 simulations.

#### 5.3 Comparison of the estimate accuracies

Once tuning parameters are obtained, we focus on the fixed–effect estimates' accuracy. Since the response–vectors  $y_1$  and  $y_2$  are normally distributed and have comparable orders of magnitude, the fixed–effect relative errors are on the same scale. Then we consider a risk–averse comparison criterion called "Mean Upper Relative Squared Error" (MURSE) defined as

MURSE 
$$(\beta_1, \beta_2) = \frac{1}{B} \sum_{b=1}^{B} \max \left\{ \frac{\left\| \widehat{\beta}_1^{(b)} - \beta_1 \right\|^2}{\left\| \beta_1 \right\|^2}, \frac{\left\| \widehat{\beta}_2^{(b)} - \beta_2 \right\|^2}{\left\| \beta_2 \right\|^2} \right\},\$$

where  $\widehat{\boldsymbol{\beta}}_{\boldsymbol{k}}^{(b)}$  is the estimate of  $\boldsymbol{\beta}_{\boldsymbol{k}}$  associated with sample *b*. The MURSE values for mixed– SCGLR, LMM–ridge and GLMM–LASSO are presented in Table 2. The LMM results obtained without regularisation are also presented. They were computed using the R package lme4 (Bates et al., 2015). In the latter case, relative errors increase dramatically with  $\tau$ . Those of ridge and LASSO increase less drastically (but increase anyway) because these methods suffer from the high correlations among the explanatory variables. Except for  $\tau = 0.1$ , mixed–SCGLR provides the most accurate fixed effect estimates. Indeed, if there are no real bundles in X ( $\tau \simeq 0$ ), searching for structures in X may lead mixed– SCGLR to be slightly less accurate. Conversely, mixed–SCGLR takes advantage of the high correlations among the explanatory variables: the stronger the structures (high  $\tau$ ), the more efficient the method.

parameter values.								
$\mathbf{L}\mathbf{M}\mathbf{M}$	CLMM_LASSO	LMM_ridgo	mixed_SCCLB					

Table 2: Mean Upper Relative Squared Error (MURSE) values associated with the optimal

	(no regularisation)	GLMM-LASSO	LMM–ridge	mixed–SCGLR
$\tau = 0.1$	0.12	0.05	0.08	0.12
$\tau = 0.3$	0.33	0.12	0.13	0.10
$\tau = 0.5$	0.61	0.20	0.16	0.07
$\tau = 0.7$	1.32	0.25	0.20	0.06
$\tau = 0.9$	4.62	0.26	0.31	0.05

#### 5.4 Model interpretation

This section aims at highlighting the power of mixed–SCGLR for model interpretation. Figure 2 presents an example of the first component planes obtained for  $\tau = 0.5$ , with associated optimal parameter values  $s^* = 0.7$  and  $K^* = 3$ . We still impose l = 4. The first two components obtained are the ones which explain the responses. It clearly appears that  $y_1$  is explained by bundle  $X_1$  and  $y_2$  by  $X_2$ . Interestingly, although bundle  $X_0$  is the one with maximum inertia (26.83%), it appears only along the third component, for having no explanatory part.



Figure 2: Component planes (1, 2) and (1, 3) given by mixed-SCGLR on simulated data. The black arrows represent the explanatory variables. The red ones represent the projection of the X-part of the linear predictors associated with  $y_1$  and  $y_2$ . The percentage of inertia captured by each component is given in parentheses.

#### 5.5 Additional simulations involving non–Gaussian outcomes

This section aims at assessing our method in the case of Bernoulli ( $\mathcal{B}$ ) and Poisson ( $\mathcal{P}$ ) distributions of responses. We still consider N = 10 groups and R = 10 observations per group. We keep design matrices X and U defined in Section 5.1, as well as the values of  $\beta_1$ ,  $\beta_2$ ,  $\sigma_1^2$  and  $\sigma_2^2$ . The group variance components are given by  $\varsigma_1^2 = 0.1\sigma_1^2$  and  $\varsigma_2^2 = \sigma_2^2$  so that for each  $k \in \{1, 2\}$ ,  $\tilde{\xi}_k \sim \mathcal{N}_N(\mathbf{0}, \varsigma_k^2 \mathbf{I}_N)$ . Then given  $\tilde{\xi}_1$  and  $\tilde{\xi}_2$ , we simulate  $Y = [\mathbf{y}_1 | \mathbf{y}_2]$  as

$$\begin{cases} \boldsymbol{y_1} \sim \mathcal{B}\left(\boldsymbol{p} = \text{logit}^{-1}\left[\boldsymbol{X}\boldsymbol{\theta_1} + \boldsymbol{U}\widetilde{\boldsymbol{\xi_1}}\right]\right) \\ \boldsymbol{y_2} \sim \mathcal{P}\left(\boldsymbol{\lambda} = \exp\left[\boldsymbol{X}\boldsymbol{\theta_2} + \boldsymbol{U}\widetilde{\boldsymbol{\xi_2}}\right]\right), \end{cases}$$
(14)

where  $\theta_1 = 0.1\beta_1$  and  $\theta_2 = \beta_2$ . Again, for each value of  $\tau$ , B = 100 samples are generated according to Model (14). As in Section 5.2, tuning parameters are calibrated so as to minimise the cross-validation error (13). However, since  $y_1$  and  $y_2$  do not have the same range of values, the prediction errors have to be standardised. The cross-validation error for response  $y_k$  in the *b*-th sample is now given by

$$E_k^{(b)} = \frac{1}{5} \sum_{j=1}^5 \sqrt{\frac{1}{20} \sum_{i \in \mathcal{P}_j} \frac{\left(y_{k,i}^{(b)} - \widehat{y_{k,i(-j)}^{(b)}}\right)^2}{\widehat{\operatorname{var}}\left(\widehat{y_{k,i(-j)}^{(b)}}\right)^2}}.$$

Unlike in Section 5.3, the response-vectors do not come from the same distribution and have different orders of magnitude. The fixed–effect relative errors are thus not comparable. To compare mixed–SCGLR with GLMM–LASSO and classical GLMM (without regularisation), we thus use the Mean Relative Squared Error (MRSE) defined as

MRSE 
$$(\boldsymbol{\theta}_{k}) = \frac{1}{B} \sum_{b=1}^{B} \frac{\left\| \widehat{\boldsymbol{\theta}}_{k}^{(b)} - \boldsymbol{\theta}_{k} \right\|^{2}}{\left\| \boldsymbol{\theta}_{k} \right\|^{2}}, \ k \in \{1, 2\},$$

where  $\hat{\theta}_{k}^{(b)}$  is the estimate of  $\theta_{k}$  from the *b*-th sample. MRSE values for the GLMM, mixed-SCGLR and GLMM-LASSO are presented in Table 3. For all methods, estimating a Bernoulli model is obviously a more challenging task than estimating a Poisson model. Regardless of the level of redundancy  $\tau$ , both mixed-SCGLR and GLMM-LASSO outperform classical GLMM estimation. Compared with the Gaussian case (Section 5.3), the results deteriorate but (overall) the same behaviours are observed.

- ► For  $\tau = 0.1$ , fixed-effect estimates provided by mixed-SCGLR are less accurate than those provided by GLMM-LASSO. In this case, GLMM-LASSO has indeed a double advantage. First, many  $\theta_{k,j}$ 's are true zeros. Unlike mixed-SCGLR, GLMM-LASSO often shrinks their estimates to exactly zero. Second, since the level of redundancy is low, GLMM-LASSO also provides accurate coefficient estimates of active variables.
- ▶ By contrast, for  $\tau \ge 0.3$ , mixed–SCGLR takes advantage of redundancies within the explanatory variables. Thus, mixed-SCGLR outperforms GLMM–LASSO in this case, despite the sparse structure of the  $\theta_k$ 's.

Even if the response variables are not Gaussian, the power of mixed–SCGLR for model interpretation is preserved. Graphical diagnoses similar to those provided in Section 5.4 are available in the Supplementary Material.

	GLN (no regula	MM arisation)	GLMM-	LASSO	mixed-SCGLR		
	Bernoulli	Poisson	Bernoulli	Poisson	Bernoulli	Poisson	
$\tau = 0.1$	316.48	0.54	8.61	0.30	14.71	0.46	
$\tau = 0.3$	398.78	0.64	9.23	0.36	7.21	0.21	
$\tau = 0.5$	576.68	0.87	14.48	0.44	2.01	0.09	
$\tau = 0.7$	886.04	1.28	17.37	0.47	1.50	0.07	
$\tau = 0.9$	2840.10	3.72	17.24	0.59	1.31	0.05	

Table 3: Mean Relative Squared Error (MRSE) values obtained with Bernoulli and Poisson responses.

### 6 An application to forest ecology data

#### 6.1 Data description

The present study is based on the *Genus* dataset of the CoForChange project (see http: //www.coforchange.eu). The subsample we consider gives the abundance of 8 common tree genera on 2615 Congo Basin land plots. These plots are grouped into 22 forest concessions. To predict abundances, we have 56 environmental variables, plus 2 explanatory variables which code geology and anthropogenic interference. X consists of all environmental variables which are:

- ▶ 29 physical factors linked to topography, rainfall or soil moisture,
- ▶ 25 photosynthesis activity indicators (the Enhanced Vegetation Indices, EVI, the Near–InfraRed indices, NIR, and the Mid–InfraRed indices, MIR),
- $\blacktriangleright$  2 indicators which describe the tree height.

Physical factors are many and redundant: monthly rainfalls are highly correlated, and so are photosynthesis activity indicators. By contrast, geology and anthropogenic interference are weakly correlated and interesting per se. These variables are then considered as additional explanatory variables and included in matrix A.

#### 6.2 Model and parameter calibration

Abundances of species given in *Genus* are count data. For each  $k \in \{1, ..., 8\}$ , we consider a Poisson regression with log link

$$oldsymbol{y_k} \sim \mathcal{P}\left(oldsymbol{\lambda} = \exp\left[\sum_{j=1}^K \left(oldsymbol{X}oldsymbol{u_j}
ight)\gamma_{k,j} + oldsymbol{A}oldsymbol{\delta_k} + oldsymbol{U}oldsymbol{\xi_k}
ight]
ight),$$

where  $\boldsymbol{\xi}_k$  is the 22-level random-effect vector used to model the dependence between the observations of  $\boldsymbol{y}_k$  within concessions. The first cross-validations we performed — with different fixed values of parameters s and l — indicated that four components were sufficient to capture most of the information in  $\boldsymbol{X}$  needed to model and predict responses. We therefore keep  $K^* = 4$ . The optimal values of trade-off and locality parameter  $s^*$  and  $l^*$  are then determined through another cross-validation. Using the same procedure and notations as in Section 5.2, the data are divided into five parts  $\mathcal{P}_1, \ldots, \mathcal{P}_5$ . Let  $n_j$  be the size of  $\mathcal{P}_j$ .

$$E = \frac{1}{8} \sum_{k=1}^{8} E_k,$$

where

$$E_k = \frac{1}{5} \sum_{j=1}^5 \sqrt{\frac{1}{n_j} \sum_{i \in \mathcal{P}_j} \frac{\left(y_{k,i} - \widehat{y_{k,i(-j)}}\right)^2}{\widehat{\operatorname{var}}\left(\widehat{y_{k,i(-j)}}\right)}}.$$
(15)

On Figure 3, we plot the errors E for parameter pairs  $(s, l) \in \mathcal{E}_s \times \mathcal{E}_l$ , where

$$\mathcal{E}_s = \{0.025, 0.1, 0.2, \dots, 1\}$$
$$\mathcal{E}_l = \{1, 2, \dots, 10, 12, 14, \dots, 30, 35, 40, 45, 50\}$$

Parameter grid  $\mathcal{E}_s \times \mathcal{E}_l$  therefore contains 264 pair values. Selecting the best parameter pair from  $\mathcal{E}_s \times \mathcal{E}_l$  through a 5-fold cross-validation requires a computation time of about 65 minutes (parallel computing on 6 CPU cores, Intel Core i7-6700HQ, 2.6GHz). It should be noted that there is a risk of non-convergence when the trade-off parameter s is too close to 0. Indeed, if we consider no structural information (s exactly equal to 0) in  $\mathbf{X}$ , mixed-SCGLR merely performs classical GLMM estimation and does not converge with this data. When s = 0.025, our algorithm converges but leads to fairly unstable estimates and high cross-validation errors because regularisation is then very weak. By contrast, the components calculated with  $s \in \{0.5, 0.6, \dots, 1\}$  are close to principal components. The associated errors are therefore stable in most cases, but rather high. Finally,  $s \in \{0.1, \dots, 0.4\}$  leads to the lowest cross-validation errors, but only for  $l \leq 10$ . Indeed, when s is not too high, mixed-SCGLR may focus on the most predictive structures of X. However, parameter l must not exceed a certain value, in order to avoid being drawn towards too local variable-bundles. As can be seen, choosing  $(s^*, l^*) = (0.1, 10)$  minimises the cross-validation error.



Figure 3: Behaviour of the cross-validation error E for trade-off parameter  $s \in \{0.025, 0.1, 0.2, \dots, 1\}$ , as a function of locality parameter  $l \in [1, 50]$ .

#### 6.3 Prediction quality and interpretation results

This part evaluates the benefits obtained by taking within-group dependence into account. The predictions we get with mixed-SCGLR and with initial version of SCGLR are compared with respect to the cross-validation criterion given by (15). Table 4 summarises the  $E_k$ 's for both SCGLR and mixed-SCGLR methods. Optimal parameter value triplet  $(K^*, s^*, l^*) =$  (4, 0.1, 10) is selected for both methods. For each  $k \in \{1, \dots, 8\}$ , mixed–SCGLR gives a lower cross–validation error than SCGLR: taking into account the within–group dependence has clearly improved prediction performances.

	$E^1_{\rm cv}$	$E_{\rm cv}^2$	$E_{\rm cv}^3$	$E_{\rm cv}^4$	$E_{\rm cv}^5$	$E_{\rm cv}^6$	$E_{\mathrm{cv}}^7$	$E_{\rm cv}^8$
SCGLR	1.32	2.46	3.27	1.43	2.56	1.28	1.54	3.44
mixed–SCGLR	1.24	1.95	2.92	1.32	2.27	1.15	1.31	3.01

Table 4: Cross–validation errors for each response variable.

Moreover, mixed–SCGLR enables to correctly reconstitute observed abundance maps, as illustrated on Figure 4.



Figure 4: Abundance maps issued from mixed–SCGLR. The plots respectively show real abundance (left) and associated conditional predictions (right) of the tree species number 8. Each point represents a land plot (2615 in total).

As has been seen in Section 5.4, mixed–SCGLR allows an easy interpretation of the model through the decomposition of linear predictors on interpretable components. Figure 5 shows the first two component planes resulting from mixed–SCGLR on real data *Genus*. Component plane (1, 2) reveals two patterns. The first one is a global rain–wind

pattern driven by the *pluvio*'s and *wd*'s variables which explain the abundances of Species 1, 2, 5, 6. The second is a rather local pattern driven by variables *altitude*, *wetness* and annual pluviometry (*pluvio\_an*) which prove important to model and predict responses  $y_3$  and  $y_7$ . Lastly, Component 3 reveals a photosynthesis pattern driven by a part of the *Evi*'s, which seems useful to predict  $y_4$  and  $y_8$ .



Figure 5: Component planes (1, 2) and (1, 3) output by mixed-SCGLR on dataset Genus, with optimal parameter triplet  $(K^*, s^*, l^*) = (4, 0.1, 10)$ . The left-hand side plot displays only variables having cosine greater than 0.7 with component plane (1, 2). The right-hand side plots variables having cosine greater than 0.75 with component plane (1, 3).

The decomposition of linear predictors on interpretable components allows to detect the species that tend to share common explanatory dimensions and those which are more idiosyncratic. We can then identify the variable-bundles these dimensions are related to. The underlying goal is a better understanding of the bio- and ecosystem diversity with a view to preserve them. Species 1, 2, 5 and 6 are sensitive to the same rain-wind regime, and Species 4 and 8 are explained by the same photosynthetic pattern. On the contrary, Species 3 and 7 are clearly separated. Species 7 grows at high altitudes where the atmosphere is rather dry while the abundance of Species 3 is favoured by regular rainfall and high humidity.

# 7 Discussion and Conclusions

Like Sufficient Dimension Reduction (SDR) methods, mixed–SCGLR is based on the construction of a reduction function of dimension less than p which tries to capture all the relevant information that X contains about Y. However, the two approaches do not exactly pursue the same objectives. Indeed, SDR methods look for the "central subspace" containing the predictive information irrespective of the structures within X (e.g. dimensions capturing a large part of X's variance, or bundles of correlated variables). Mixed–SCGLR rather aims at basing the explanatory subspace on such structural dimensions so as to both gain interpretability and stabilise prediction. We think that extracting a hierarchy of strong and interpretable dimensions, and decomposing the linear predictor on them, is an essential asset in model–building. The difference in goals entails a difference in means: SDR is based on the sufficiency principle, which is enough to identify a subspace but not to track strong predictive dimensions in it. By contrast, in the wake of PLS regression, mixed–SCGLR uses a criterion combining goodness–of–fit and structural relevance of components.

The supervised–component paradigm has proved effective in situations where regularisation is necessary but where variable selection is inappropriate — for instance when the true explanatory dimensions are latent and indirectly measured through highly correlated proxies.

- ▶ When l = 1, trade-off parameter s allows to continuously tune the attraction of components towards the principal components of explanatory variables. This results in a continuum between classical GLMM estimation (s = 0 is associated with no regularisation) and principal component generalised linear mixed regression (with s = 1).
- ▶ When l > 1, we take better advantage of local predictive structures in X. The components we build are then usually closer to local gatherings of variables, thus easier to interpret.

Mixed–SCGLR is able to identify more or less local predictive structures common to all the  $y_k$ 's and performs well on grouped data with Gaussian, Bernoulli, binomial and Poisson

outcomes. Compared to penalty-based approaches as ridge or LASSO, the orthogonal components built by mixed–SCGLR reveal the multidimensional explanatory and predictive dimensions, and greatly facilitate the interpretation of the model.

However, a natural question arises as to the accuracy of our methodology under significant deviations from normality. With binary data for instance, variance component estimates are prone to some bias towards zero (McCulloch, 1997). That is why other estimation strategies might be considered, especially Monte Carlo integration methods which have the advantage of being based on direct approximations of the likelihood. Some examples are the MCMC methods developed by Hadfield (2010) in the GLMM framework, and the Monte Carlo Likelihood Approximation (MCLA) proposed by Knudson (2016). Indirect maximisations of the likelihood are also available such as Monte Carlo Expectation– Maximisation (MCEM) and Monte Carlo Newton–Raphson (McCulloch, 1997). We think that these methodologies and Schall's could be combined sequentially. Indeed we could first take advantage of the linear approximation of the model in order to build the components, and then use MC–based methods to estimate both fixed–effect parameters and variance components. This would lead to replacing the current iteration of mixed–SCGLR given by Algorithm 1 with the following steps (to keep things simple, we take the canonical link):

- 1. Compute components  $F = [f_1 | \dots | f_K]$  via the PING algorithm on Schall's linearised models.
- 2. For each  $k \in \{1, \ldots, q\}$ , consider the hierarchy

$$\begin{split} f_{\boldsymbol{y_k}|\boldsymbol{\xi_k},\boldsymbol{\gamma_k},\boldsymbol{\delta_k}}\left(\boldsymbol{y_k}|\boldsymbol{\xi_k},\boldsymbol{\gamma_k},\boldsymbol{\delta_k}\right) &= \exp\left\{\boldsymbol{y_k^{\mathrm{T}}\boldsymbol{\eta_k^{\xi}}} - \boldsymbol{1^{\mathrm{T}}c}\left(\boldsymbol{\eta_k^{\xi}}\right) + \boldsymbol{1^{\mathrm{T}}d}\left(\boldsymbol{y_k}\right)\right\} \\ \boldsymbol{\xi_k}|\boldsymbol{D_k} \sim \mathcal{N}\left(\boldsymbol{0},\boldsymbol{D_k}\right), \end{split}$$

where  $\eta_k^{\boldsymbol{\xi}} = \boldsymbol{F} \boldsymbol{\gamma}_k + \boldsymbol{A} \boldsymbol{\delta}_k + \boldsymbol{U} \boldsymbol{\xi}_k$ , and c, d are the functions associated with the natural parametrisation of the GLM. For example, for the Bernoulli–logistic regression, we have:  $c(x) = \log(1 + e^x)$  and d(x) = 0.

- 3. Apply MC-based methods such as MCMC, MCLA, MCEM or MCNR to update  $\gamma_k$ ,  $\delta_k, \xi_k$  and  $D_k, k \in \{1, \dots, q\}$ .
- 4. Update working variables and weight matrices to define the new Schall's linearised models.

Even though such MC–based methods are computationally much more intensive than the "Joint–Maximisation" and have intrinsic disadvantages (particularly in the assessment of convergence and in the choice of prior distributions), they could give better results in case of binary data.

#### SUPPLEMENTARY MATERIAL

- Additional simulations: The first simulation reproduces that of Section 5.5 in the case of binomial and Poisson outcomes. The second simulation explores a different structure of variable–bundles, considers Gaussian, binomial and Poisson outcomes, and presents results concerning variance component estimates. The third one involves high dimensional data. (pdf file)
- **Projected Iterated Normed Gradient (PING) algorithm:** We give some technical details about the PING algorithm, which maximises, at least locally, any criterion on the unit sphere. (pdf file)
- **R package mixedSCGLR:** We provide an R package to perform mixed-SCGLR, also available at https://github.com/SCnext/mixedSCGLR. It contains the dataset *Genus* used in Section 6. The package also provides demo codes, in particular for visualising the component planes (mixedSCGLR.tar.gz).
- **Code for running simulations:** We also provide the R codes required to reproduce most of the simulation results (R and Rdata files).

#### ACKNOWLEDGMENTS

The extended data *Genus* required the arrangement and the inventory of 140.000 developed plots across four countries : Central African Republic, Gabon, Cameroon and Democratic Republic of Congo. The authors thank the members of the CoForTips project for allowing the use of this data. We are also grateful to the editor, the associate editor and to the referees for their thorough and constructive review of this work.

# Supplementary file to "Component-based regularisation of multivariate generalised linear mixed models": THE PING ALGORITHM

Jocelyn Chauvet<sup>1</sup> Catherine Trottier<sup>1,2</sup> Xavier  $Bry^1$ 

<sup>1</sup> IMAG, Univ Montpellier, CNRS, Montpellier, France. jocelyn.chauvet@umontpellier.fr ; xavier.bry@umontpellier.fr <sup>2</sup> Univ Paul-Valéry Montpellier 3, Montpellier, France. catherine.trottier@univ-montp3.fr

The Projected Iterated Normed Gradient (PING) is a basic extension of the iterated power algorithm, for solving any program of the form

$$\begin{cases} \max \quad \mathcal{J}_{h}\left(\boldsymbol{u}\right), \\ \text{subject to:} \quad \boldsymbol{u}^{\mathrm{T}}\boldsymbol{M}^{-1}\boldsymbol{u} = 1 \text{ and } \boldsymbol{\Delta}_{h}^{\mathrm{T}}\boldsymbol{u} = \boldsymbol{0}. \end{cases}$$
(16)

Note that putting  $\boldsymbol{v} = \boldsymbol{M}^{-1/2}\boldsymbol{u}$ ,  $\mathcal{G}_h(\boldsymbol{v}) = \mathcal{J}_h\left(\boldsymbol{M}^{1/2}\boldsymbol{v}\right)$  and  $\boldsymbol{B}_h = \boldsymbol{M}^{1/2}\boldsymbol{\Delta}_h$ , Program (16) is strictly equivalent to Program (17):

$$\begin{cases} \max \quad \mathcal{G}_{h}\left(\boldsymbol{v}\right), \\ \text{subject to:} \quad \boldsymbol{v}^{\mathrm{T}}\boldsymbol{v} = 1 \text{ and } \boldsymbol{B}_{h}^{\mathrm{T}}\boldsymbol{v} = \boldsymbol{0}. \end{cases}$$
(17)

Denoting

$$\Pi_{\text{span}\{\boldsymbol{B}_{\boldsymbol{h}}\}^{\perp}} = \boldsymbol{I} - \boldsymbol{B}_{\boldsymbol{h}} \left(\boldsymbol{B}_{\boldsymbol{h}}^{\text{T}} \boldsymbol{B}_{\boldsymbol{h}}\right)^{-1} \boldsymbol{B}_{\boldsymbol{h}}^{\text{T}} \text{ and}$$
$$\Gamma_{\boldsymbol{h}} \left(\boldsymbol{v}\right) = \bigvee_{\boldsymbol{v}} \mathcal{G}_{\boldsymbol{h}} \left(\boldsymbol{v}\right),$$

a Lagrange multiplier-based reasoning gives the basic iteration of the PING algorithm:

$$\boldsymbol{v}^{[t+1]} = \frac{\Pi_{\operatorname{span}\{\boldsymbol{B}_{\boldsymbol{h}}\}^{\perp}} \Gamma_{h}\left(\boldsymbol{v}^{[t]}\right)}{\left\|\Pi_{\operatorname{span}\{\boldsymbol{B}_{\boldsymbol{h}}\}^{\perp}} \Gamma_{h}\left(\boldsymbol{v}^{[t]}\right)\right\|}.$$
(18)

Although Iteration (18) follows a direction of ascent, it does not guarantee that  $\mathcal{G}_h$  actually increases on every step. We therefore propose a generic iteration of PING (Algorithm 2) and an alternative one (Algorithm 3), which both ensure that the criterion increases.

while convergence of 
$$\boldsymbol{v}$$
 non reached do  

$$\boldsymbol{\kappa}^{[t]} = \frac{\Pi_{\text{span}\{\boldsymbol{B}_{h}\}^{\perp}} \Gamma_{h}\left(\boldsymbol{v}^{[t]}\right)}{\left\|\Pi_{\text{span}\{\boldsymbol{B}_{h}\}^{\perp}} \Gamma_{h}\left(\boldsymbol{v}^{[t]}\right)\right\|}$$
A unidimensional Newton–Raphson maximisation procedure is used to find  
the maximum of  $\mathcal{G}_{h}\left(\boldsymbol{v}\right)$  on the arc  $\left(\boldsymbol{v}^{[t]}, \boldsymbol{\kappa}^{[t]}\right)$  and take it as  $\boldsymbol{v}^{[t+1]}$ .  
 $t \leftarrow t+1$   
end  
Algorithm 2: Generic iteration of the PING algorithm

while convergence of 
$$\boldsymbol{v}$$
 non reached do  

$$\begin{array}{c|c} \boldsymbol{m} \leftarrow & \frac{\Pi_{\operatorname{span} \{\boldsymbol{B}_{h}\}^{\perp}} \Gamma_{h} \left( \boldsymbol{v}^{[t]} \right)}{\left\| \Pi_{\operatorname{span} \{\boldsymbol{B}_{h}\}^{\perp}} \Gamma_{h} \left( \boldsymbol{v}^{[t]} \right) \right\|} \\ \text{while } \mathcal{G}_{h} \left( \boldsymbol{m} \right) < \mathcal{G}_{h} \left( \boldsymbol{v}^{[t]} \right) \operatorname{do} \\ \left\| \boldsymbol{m} \leftarrow & \frac{\boldsymbol{v}^{[t]} + \boldsymbol{m}}{\left\| \boldsymbol{v}^{[t]} + \boldsymbol{m} \right\|} \\ \text{end} \\ \boldsymbol{v}^{[t+1]} = \boldsymbol{m} \\ t \leftarrow t+1 \\ \text{end} \end{array}$$
Algorithm 3: Alternative generic iteration of the PING algorithm

First rank component. Component  $f_1 = X u_1$  is obtained by solving

$$\begin{cases} \max \quad s \log \left[ \phi \left( \boldsymbol{u} \right) \right] + (1 - s) \log \left[ \psi_{\boldsymbol{A}} \left( \boldsymbol{u} \right) \right] \\ \text{subject to:} \quad \boldsymbol{u}^{\mathrm{T}} \boldsymbol{M}^{-1} \boldsymbol{u} = 1. \end{cases}$$

This corresponds to Program (16) with h = 0, where

- $\blacktriangleright \mathcal{J}_0(\boldsymbol{u}) = s \log \left[\phi(\boldsymbol{u})\right] + (1-s) \log \left[\psi_{\boldsymbol{A_0}}(\boldsymbol{u})\right],$
- $\blacktriangleright$   $A_0 = A$  (the matrix of additional explanatory variables), and
- $\blacktriangleright \ \Delta_0 = 0.$

In this particular case, we have  $\boldsymbol{B_0} = \boldsymbol{M}^{1/2} \boldsymbol{\Delta_0} = \boldsymbol{0},$  and so:

$$\Pi_{\operatorname{span}\{B_0\}^{\perp}} = I$$

Higher rank components. Let  $F_h = \begin{bmatrix} f_1 & \dots & f_h \end{bmatrix}$  be the matrix of the first h components and  $A_h = \begin{bmatrix} F_h & A \end{bmatrix}$ . Let P denote the weight matrix reflecting the a priori relative importance of observations ( $P = \frac{1}{n}I_n$  if all observations are of equal importance). Component  $f_{h+1} = Xu_{h+1}$  is obtained by solving

$$\begin{cases} \max \quad s \log \left[ \phi \left( \boldsymbol{u} \right) \right] + (1 - s) \log \left[ \psi_{\boldsymbol{A_h}} \left( \boldsymbol{u} \right) \right] \\ \text{subject to:} \quad \boldsymbol{u}^{\mathrm{T}} \boldsymbol{M}^{-1} \boldsymbol{u} = 1 \text{ and } \boldsymbol{F_h^{\mathrm{T}}} \boldsymbol{P} \boldsymbol{X} \boldsymbol{u} = \boldsymbol{0} \end{cases}$$

This corresponds to Program (16), where

►  $\mathcal{J}_h(\boldsymbol{u}) = s \log [\phi(\boldsymbol{u})] + (1-s) \log [\psi_{\boldsymbol{A}_h}(\boldsymbol{u})],$ ►  $\boldsymbol{A}_h = [\boldsymbol{F}_h \mid \boldsymbol{A}], \text{ and}$ ►  $\boldsymbol{\Delta}_h = \boldsymbol{X}^{\mathrm{T}} \boldsymbol{P} \boldsymbol{F}_h.$ 

**Initialisation.** To quickly find  $f_1$ , algorithm PING is initialised with the first PLS component of the responses on X. In like manner, for  $h \ge 2$ , PING is initialised with the first PLS component of the responses on X deflated on components  $F_{h-1} = [f_1 | \dots | f_{h-1}]$ .

# Supplementary file to "Component-based regularisation of multivariate generalised linear mixed models": ADDITIONAL SIMULATIONS

Jocelyn Chauvet<sup>1</sup> Catherine Trottier<sup>1,2</sup> Xavier  $Bry^1$ 

<sup>1</sup> IMAG, Univ Montpellier, CNRS, Montpellier, France. jocelyn.chauvet@umontpellier.fr; xavier.bry@umontpellier.fr <sup>2</sup> Univ Paul-Valéry Montpellier 3, Montpellier, France. catherine.trottier@univ-montp3.fr

# 8 Comparative results with binomial and Poisson out-

#### comes

In this section, we simply extend the simulation scheme presented in Section 5.5 to binomial and Poisson outcomes. We maintain design matrices X and U as defined in Section 5.1. Fixed-effect parameters  $\theta_k$ 's and random-effect vectors  $\tilde{\xi}_k$ 's are defined in Section 5.5. Given  $\tilde{\xi}_1$  and  $\tilde{\xi}_2$ , we then simulate  $Y = [y_1 | y_2]$  as

$$\begin{cases} \boldsymbol{y_1} \sim \mathcal{B}in\left(\text{trials} = 50\,\boldsymbol{1}_n,\,\boldsymbol{p} = \text{logit}^{-1}\left[\boldsymbol{X}\boldsymbol{\theta_1} + \boldsymbol{U}\widetilde{\boldsymbol{\xi_1}}\right]\right) \\ \boldsymbol{y_2} \sim \mathcal{P}\left(\boldsymbol{\lambda} = \exp\left[\boldsymbol{X}\boldsymbol{\theta_2} + \boldsymbol{U}\widetilde{\boldsymbol{\xi_2}}\right]\right). \end{cases}$$

Table 5 gives the Mean Relative Squared Error (MRSE) values for  $\theta_1$  and  $\theta_2$  obtained on 100 samples for each value of  $\tau$ .

For the Poisson distribution, the results in Table 5 are essentially identical to those in the article: mixed-SCGLR outperforms GLMM-LASSO (Groll, 2017, R package glmmLasso) except for  $\tau = 0.1$ . As for the binomial distribution, the regularisation provided by mixed-SCGLR improves the results obtained without regularisation (Bates et al., 2015, R package lme4), regardless of the level of redundancy within the explanatory variables. Unsurprisingly, the errors are much smaller than in the binary case.

Table 5: Mean Relative Squared Error (MRSE) values obtained with binomial and Poisson distributions. The R package glmmLasso does not handle binomial outcomes but only Bernoulli ones, which precludes comparison in this case.

$\operatorname{GLMM}$			CIMM IASSO	mined SCCLD		
	(no regularisation)		GLIVINI–LASSO	IIIIxeu–50GLK		
	Binomial	Poisson	Poisson	Binomial	Poisson	
$\tau = 0.1$	2.31	0.50	0.31	0.51	0.45	
$\tau = 0.3$	3.07	0.60	0.33	0.28	0.18	
$\tau = 0.5$	3.93	0.75	0.39	0.15	0.09	
$\tau = 0.7$	6.50	1.07	0.40	0.10	0.07	
$\tau = 0.9$	19.29	2.71	0.42	0.07	0.05	

The power of mixed–SCGLR in terms of model interpretation remains the same for non–Gaussian outcomes. Figure 6 (respectively Figure 7) presents an example of the first component planes output by mixed–SCGLR in the binomial/Poisson (respectively Bernoulli/Poisson) case. As for Gaussian outcomes, the component planes reveal that  $y_1$ is explained by bundle  $X_1$  and  $y_2$  by  $X_2$ . In the binomial/Poisson case with  $\tau = 0.3$ (Figure 6), predictive bundles  $X_1$  and  $X_2$  are captured respectively by the first and the second components. The third component aligns on nuisance bundle  $X_0$ , despite its high inertia. Figure 7 illustrates what may happen when the level of redundancy is very high ( $\tau = 0.7$  here). Since the explanatory variables are highly correlated, mixed–SCGLR regularisation requires that the structural relevance be given a heavy weight with respect to the goodness–of–fit, which leads to a trade–off parameter s close to 1 (s = 0.9 here). Having the greatest structural strength, the nuisance bundle is captured by the second component despite its lack of explanatory power. This is sometimes the price to be paid for the trade– off. In our example, the second explanatory bundle is captured by the third component, so that the predictive dimensions are accurately represented in component plane (1, 3).



Figure 6: Example of component planes given by mixed–SCGLR in the binomial/Poisson case for  $\tau = 0.3$ , with parameter triplet (K, s, l) = (3, 0.5, 2).



Figure 7: Example of component planes given by mixed–SCGLR in the Bernoulli/Poisson case for  $\tau = 0.7$ , with parameter triplet (K, s, l) = (3, 0.9, 4).

# 9 A new structure for the bundles — Gaussian, Poisson and binomial distributions

This simulation study tests mixed-SCGLR on a slightly more complex bundle structure. Results concerning variance component estimates are also presented.

We consider a fixed-effect design matrix  $X_{n \times p}$  partitioned into 3 blocks  $\mathcal{G}_1$ ,  $\mathcal{G}_2$  and  $\mathcal{G}_3$ . Block  $\mathcal{G}_1$  contains 10 predictive explanatory variables structured about a latent variable  $\varphi_1 \sim \mathcal{N}_n(\mathbf{0}, \sigma_{\text{LV}}^2 \mathbf{I}_n)$ . Thus for each  $j \in \{1, \ldots, 10\}$ ,  $x_j = \varphi_1 + \varepsilon_j$ , where  $\varepsilon_j \sim \mathcal{N}_n(\mathbf{0}, \sigma_{\text{noise}}^2 \mathbf{I}_n)$  such that  $\sigma_{\text{LV}}^2 + \sigma_{\text{noise}}^2 = 1$ . The correlation within  $\mathcal{G}_1$  is tuned by signal to noise (StN) ratio  $\sigma_{\text{LV}}^2/\sigma_{\text{noise}}^2$  (chosen in  $\{\frac{1}{3}, 1, 3\}$  in practice).  $\mathcal{G}_2$  contains a single predictive variable  $\varphi_2 = x_{11} \sim \mathcal{N}_n(\mathbf{0}, \mathbf{I}_n)$ . For each  $k \in \{1, 2, 3\}$ , random-effect vectors are simulated as  $\boldsymbol{\xi}_k \stackrel{\text{ind.}}{\sim} \mathcal{N}_N(\mathbf{0}, \sigma_k^2 \mathbf{I}_N)$ . Given  $\boldsymbol{\xi}_1, \boldsymbol{\xi}_2, \boldsymbol{\xi}_3$ , we simulate 3 responses having different distributions,  $\boldsymbol{Y} = [\boldsymbol{y}_1 | \boldsymbol{y}_2 | \boldsymbol{y}_3]$ , as

$$\begin{cases} \boldsymbol{y_1} \sim \mathcal{N}_n \Big( \boldsymbol{\mu} = \alpha_1 \boldsymbol{\varphi_1} + \boldsymbol{U} \boldsymbol{\xi_1}, \, \boldsymbol{\Sigma} = \boldsymbol{I}_n \Big) \\ \boldsymbol{y_2} \sim \mathcal{P} \Big( \boldsymbol{\lambda} = \exp \left[ \alpha_2 \boldsymbol{\varphi_2} + \boldsymbol{U} \boldsymbol{\xi_2} \right] \Big) \\ \boldsymbol{y_3} \sim \mathcal{B}in \Big( \mathbf{trials} = 25 \, \boldsymbol{1}_n, \, \boldsymbol{p} = \mathrm{logit}^{-1} \Big[ \alpha_3 \left( \boldsymbol{\varphi_1} + \boldsymbol{\varphi_2} \right) + \boldsymbol{U} \boldsymbol{\xi_3} \Big] \Big). \end{cases}$$

In our simulations, we set  $\alpha_1 = \sigma_1^2 = 2$ ,  $\alpha_2 = \sigma_2^2 = 1$  and  $\alpha_3 = \sigma_3^2 = 0.5$ .

We consider in turn N = 10 and N = 50 groups, and R = 10 observations per group (n = 100 and n = 500 observations in total). B = 100 samples are generated for each pair of values (N, StN). The main goal of the study is to assess the ability of mixed–SCGLR to track down both latent variable  $\varphi_1$  and predictive variable  $\varphi_2$ . For j = 1 and 2, we then define

$$\operatorname{cor}_{j} = \frac{1}{B} \sum_{b=1}^{B} \left| \operatorname{cor} \left( \boldsymbol{\varphi}_{j}, \boldsymbol{f}_{j}^{(b)} \right) \right|,$$

where  $f_{j}^{(b)}$  is the component most correlated with  $\varphi_{j}$  issued from mixed–SCGLR in the *b*-th sample. Consistency of fixed–effect estimates is assessed through criteria err<sub>1</sub>, err<sub>2</sub> and  $err_3$  defined by

$$\operatorname{err}_{j} = \frac{1}{B} \sum_{b=1}^{B} \frac{\left\| \alpha_{j} \boldsymbol{\varphi}_{j} - \boldsymbol{X} \widehat{\boldsymbol{\beta}}_{j}^{(b)} \right\|^{2}}{\left\| \alpha_{j} \boldsymbol{\varphi}_{j} \right\|^{2}}, \ j \in \{1, 2\}$$
$$\operatorname{err}_{3} = \frac{1}{B} \sum_{b=1}^{B} \frac{\left\| \alpha_{3} \left( \boldsymbol{\varphi}_{1} + \boldsymbol{\varphi}_{2} \right) - \boldsymbol{X} \widehat{\boldsymbol{\beta}}_{3}^{(b)} \right\|^{2}}{\left\| \alpha_{3} \left( \boldsymbol{\varphi}_{1} + \boldsymbol{\varphi}_{2} \right) \right\|^{2}},$$

where  $\widehat{\boldsymbol{\beta}}_{\boldsymbol{j}}^{(b)}$  is the fixed-effect estimate related to response  $\boldsymbol{y}_{\boldsymbol{j}}$  associated with sample b.

Table 6 summarises the values of the afore-defined criteria and presents biases and standard errors of variance components estimates. For a given value of N, cor<sub>1</sub> increases towards 1 with ratio  $\sigma_{LV}^2/\sigma_{noise}^2$ : the tighter the block  $\mathcal{G}_1$  is structured about its latent variable, the better mixed–SCGLR can reconstruct it. The associated criterion err<sub>1</sub> then naturally decreases towards 0. On the other hand, cor<sub>2</sub> and err<sub>2</sub> are very stable, which proves that mixed–SCGLR is able to detect an isolated predictive variable among a large number of irrelevant others. As err<sub>3</sub> depends on how accurately mixed–SCGLR recovers  $\varphi_1$  and  $\varphi_2$ , it slightly decreases when the StN ratio increases. Both variance component biases and standard errors seem rather stable regardless of the value of StN. Finally, when N increases, all the cor<sub>j</sub>'s increase towards 1 and all the err<sub>j</sub>'s decrease towards 0. As far as variance component estimates are concerned, the biases are getting slightly closer to 0 and the standard errors decrease significantly.

Model interpretation is revealed by Figure 8 in the case of N = 10 groups and R = 10 observations per group. The first component aligns with block  $\mathcal{G}_1$  which alone explains response  $y_1$ . The second aligns with  $\mathcal{G}_2$  (containing single explanatory variable  $x_{11}$ ) which alone explains  $y_2$ . Finally, note that the projection of the X-part of the linear predictor related to  $y_3$  is well represented on component plane (1,2). This indicates that  $y_3$  is explained jointly by  $\mathcal{G}_1$  and  $\mathcal{G}_2$ .

	$N = 10, R = 10 \ (n = 100)$			$N = 50, R = 10 \ (n = 500)$			
$\sigma_{\rm LV}^2/\sigma_{\rm noise}^2$	$\frac{1}{3}$	1	3	$\frac{1}{3}$	1	3	
$\operatorname{cor}_1$	0.71	0.91	0.96	0.75	0.92	0.96	
$cor_2$	0.93	0.94	0.94	0.97	0.98	0.98	
$\operatorname{err}_1$	0.47	0.15	0.06	0.38	0.13	0.05	
$\operatorname{err}_2$	0.12	0.12	0.12	0.05	0.04	0.04	
$\operatorname{err}_3$	0.19	0.14	0.11	0.11	0.07	0.04	
bias $\left(\widehat{\sigma_1^2}\right)$	-0.02	-0.01	0.00	0.02	0.00	-0.02	
$\operatorname{sd}\left(\widehat{\sigma_{1}^{2}}\right)$	1.04	1.05	1.06	0.41	0.40	0.39	
bias $\left(\hat{\sigma}_2^2\right)$	-0.11	-0.08	-0.06	-0.06	-0.06	-0.06	
$\operatorname{sd}\left(\widehat{\sigma_{2}^{2}}\right)$	0.50	0.51	0.52	0.21	0.21	0.21	
bias $\left(\hat{\sigma}_3^2\right)$	-0.03	-0.04	-0.04	-0.02	-0.02	-0.02	
$\operatorname{sd}\left(\widehat{\sigma_{3}^{2}}\right)^{\prime}$	0.22	0.21	0.21	0.11	0.11	0.11	

Table 6: Summary of  $cor_j$  and  $err_j$  values, and presentation of biases and standard errors of estimated variance components.



Figure 8: Examples of the first two-component planes given by mixed-SCGLR when  $\sigma_{\text{LV}}^2/\sigma_{\text{noise}}^2 = 1/3$  (top left),  $\sigma_{\text{LV}}^2/\sigma_{\text{noise}}^2 = 1$  (top right), and  $\sigma_{\text{LV}}^2/\sigma_{\text{noise}}^2 = 3$  (bottom). When StN ratio = 1/3 (resp. StN ratio  $\in \{1, 3\}$ ), only the variables having cosine greater than 0.4 (resp. 0.5) with component plane (1, 2) are represented.

# 10 High dimensional data

#### 10.1 Key idea

To cope with high dimensional data, the key idea is to replace the fixed-effect design matrix,  $\boldsymbol{X}$ , with the matrix  $\boldsymbol{C}$  of its principal components associated with non-zero eigenvalues. More precisely,  $\lambda_j$  being the eigenvalue associated with the *j*-th eigenvector  $\boldsymbol{v}_j$ , the last eigenvector we consider,  $\boldsymbol{v}_r$ , is such that

$$\frac{\lambda_r}{\sum_{j=1}^r \lambda_j} > \frac{1}{p},$$

where p is the number of columns of matrix X. The matrix of the corresponding uniteigenvectors is denoted  $V = [v_1 | \ldots | v_r]$ , and C = XV. The component f is then sought as a combination of the principal components:  $f = Cu = X\tilde{u}$ , where  $\tilde{u} = Vu$ . Mixed-SCGLR then solves

$$\begin{cases} \max \quad s \log \left[\phi\left(\boldsymbol{u}\right)\right] + (1-s) \log \left[\psi_{\boldsymbol{A}}\left(\boldsymbol{u}\right)\right] \\ \text{subject to} \quad \boldsymbol{u}^{\mathrm{T}} \boldsymbol{C}^{\mathrm{T}} \boldsymbol{P} \boldsymbol{C} \boldsymbol{u} = 1, \end{cases}$$

where the goodness–of–fit measure,  $\psi_{A}$ , is given by

$$\psi_{\boldsymbol{A}}(\boldsymbol{u}) = \sum_{k=1}^{q} \left\| \boldsymbol{z}_{\boldsymbol{k}} \right\|_{\boldsymbol{W}_{\boldsymbol{k}}^{\boldsymbol{\xi}}}^{2} \cos^{2}_{\boldsymbol{W}_{\boldsymbol{k}}^{\boldsymbol{\xi}}} \left( \boldsymbol{z}_{\boldsymbol{k}}, \operatorname{span}\left\{ \boldsymbol{C}\boldsymbol{u}, \boldsymbol{A} \right\} \right)$$
$$= \sum_{k=1}^{q} \left\| \boldsymbol{z}_{\boldsymbol{k}} \right\|_{\boldsymbol{W}_{\boldsymbol{k}}^{\boldsymbol{\xi}}}^{2} \cos^{2}_{\boldsymbol{W}_{\boldsymbol{k}}^{\boldsymbol{\xi}}} \left( \boldsymbol{z}_{\boldsymbol{k}}, \Pi_{\operatorname{span}\left\{ \boldsymbol{C}\boldsymbol{u}, \boldsymbol{A} \right\}}^{\boldsymbol{W}_{\boldsymbol{k}}^{\boldsymbol{\xi}}} \boldsymbol{z}_{\boldsymbol{k}} \right),$$

and the structural relevance by

$$\phi\left(\boldsymbol{u}
ight) = \left[\sum_{j=1}^{p} \omega_{j}\left(\left\langle \left. \boldsymbol{C}\boldsymbol{u} \left| \left. \boldsymbol{x^{j}} \right. \right\rangle_{\boldsymbol{P}}^{2} 
ight)^{l} 
ight]^{rac{1}{l}} = \left[\sum_{j=1}^{p} \omega_{j}\left( \boldsymbol{u}^{\mathrm{T}} \left. \boldsymbol{C}^{\mathrm{T}} \boldsymbol{P} \boldsymbol{x_{j}} \boldsymbol{x_{j}}^{\mathrm{T}} \boldsymbol{P} \boldsymbol{C} \left. \boldsymbol{u} 
ight)^{l} 
ight]^{rac{1}{l}}$$

This idea is tested on simulated data where the number of explanatory variables p exceeds the number of observations n.

#### 10.2 Data generation

To generate grouped data, we consider N = 10 groups and R = 10 observations per group (i.e. a total of n = 100 observations). The random effects' design matrix is then U =  $I_N \otimes \mathbf{1}_R$ . Explanatory variables consist of four independent bundles  $X_j, j \in \{1, 2, 3, 4\}$ , such as  $\mathbf{X} = \begin{bmatrix} \mathbf{X_0} \mid \mathbf{X_1} \mid \mathbf{X_2} \mid \mathbf{X_3} \end{bmatrix}$ . Each explanatory variable is normally simulated with mean 0 and variance 1. Parameter  $\tau \in \{0.3, 0.5, 0.7\}$  tunes the level of redundancy within each bundle: the correlation matrix of bundle  $X_j$  is

$$\operatorname{cor}\left(\boldsymbol{X}_{\boldsymbol{j}}\right) = \tau \mathbf{1}_{p_{\boldsymbol{j}}} \mathbf{1}_{p_{\boldsymbol{j}}}^{\mathrm{T}} + (1-\tau) \boldsymbol{I}_{p_{\boldsymbol{j}}},$$

where  $p_j$  is the number of variables in  $X_j$ . For each  $k \in \{1, 2, 3, 4\}$ , random–effect vectors are simulated as  $\boldsymbol{\xi_k} \stackrel{\text{ind.}}{\sim} \mathcal{N}_N(\mathbf{0}, \sigma_k^2 \boldsymbol{I}_N)$ . Given  $\boldsymbol{\xi_1}, \boldsymbol{\xi_2}, \boldsymbol{\xi_3}, \boldsymbol{\xi_4}$ , we simulate 4 responses having different distributions,  $\boldsymbol{Y} = [\boldsymbol{y_1} | \boldsymbol{y_2} | \boldsymbol{y_3} | \boldsymbol{y_4}]$ , as

$$\begin{cases} \boldsymbol{y_1} \sim \mathcal{N}_n \Big( \boldsymbol{\mu} = \boldsymbol{X} \boldsymbol{\beta_1} + \boldsymbol{U} \boldsymbol{\xi_1}, \, \boldsymbol{\Sigma} = \boldsymbol{I}_n \Big) \\ \boldsymbol{y_3} \sim \mathcal{B} \Big( \boldsymbol{p} = \text{logit}^{-1} \Big[ \boldsymbol{X} \boldsymbol{\beta_2} + \boldsymbol{U} \boldsymbol{\xi_2} \Big] \Big) \\ \boldsymbol{y_3} \sim \mathcal{B} \text{in} \Big( \text{trials} = 30 \, \boldsymbol{1}_n, \, \boldsymbol{p} = \text{logit}^{-1} \Big[ \boldsymbol{X} \boldsymbol{\beta_3} + \boldsymbol{U} \boldsymbol{\xi_3} \Big] \Big) \\ \boldsymbol{y_4} \sim \mathcal{P} \Big( \boldsymbol{\lambda} = \exp \Big[ \boldsymbol{X} \boldsymbol{\beta_4} + \boldsymbol{U} \boldsymbol{\xi_4} \Big] \Big). \end{cases}$$
(19)

Response  $y_1$  is predicted only by bundle  $X_1$ ,  $y_2$  only by bundle  $X_2$ ,  $y_3$  only by bundle  $X_3$ ,  $y_4$  by both bundles  $X_2$  and  $X_3$ , and bundle  $X_0$  plays no explanatory role. Our choice for the fixed-effect parameters is

$$\beta_{1} = \left(\underbrace{0, \dots, 0}_{p_{0} \text{ times}}, \underbrace{0, 1, \dots, 0}_{p_{1} \text{ times}}, \underbrace{0, \dots, 0}_{p_{2} \text{ times}}, \underbrace{0, \dots, 0}_{p_{2} \text{ times}}, \underbrace{0, \dots, 0}_{p_{3} \text{ times}}\right)^{\mathrm{T}},$$
  
$$\beta_{2} = \left(\underbrace{0, \dots, 0}_{p_{0} \text{ times}}, \underbrace{0, \dots, 0}_{p_{1} \text{ times}}, \underbrace{0, 1, \dots, 0}_{p_{2} \text{ times}}, \underbrace{0, \dots, 0}_{p_{3} \text{ times}}, \underbrace{0, \dots, 0}_{p_{3} \text{ times}}\right)^{\mathrm{T}},$$
  
$$\beta_{3} = \left(\underbrace{0, \dots, 0}_{p_{0} \text{ times}}, \underbrace{0, \dots, 0}_{p_{1} \text{ times}}, \underbrace{0, \dots, 0}_{p_{2} \text{ times}}, \underbrace{0, 0.05, \dots, 0.05}_{p_{3} \text{ times}}\right)^{\mathrm{T}},$$
  
$$\beta_{4} = \left(\underbrace{0, \dots, 0}_{p_{0} \text{ times}}, \underbrace{0, 0.25, \dots, 0.025}_{p_{1} \text{ times}}, \underbrace{0, 0.25, \dots, 0.025}_{p_{2} \text{ times}}, \underbrace{0, \dots, 0}_{p_{3} \text{ times}}\right)^{\mathrm{T}}$$

We consider in turn p = 150 ( $p_0 = 60$ ,  $p_1 = 45$ ,  $p_2 = 30$ ,  $p_3 = 15$ ) and p = 200 ( $p_0 = 80$ ,  $p_1 = 60$ ,  $p_2 = 40$ ,  $p_3 = 20$ ) explanatory variables. Variance components are set to  $\sigma_1^2 = \sigma_2^2 = \sigma_3^2 = 0.1$ , and  $\sigma_4^2 = 0.05$ . For each value of p and for each value of  $\tau$ , B = 20 samples are generated according to Model (19).

#### 10.3 Results

Table 7 and Table 8 present the results for respectively 150 and 200 explanatory variables. They give the Mean Relative Squared Error (MRSE) values for  $\beta_k, k \in \{1, ..., 4\}$ , as well as biases and standard errors of estimated variance components, obtained on 20 samples for each value of  $\tau$ .

Table 7: Mean Relative Squared Error (MRSE) values for fixed–effect estimates, and biases and standard errors for estimated variance components, obtained with 100 observations and 150 explanatory variables.

	$eta_1$	$eta_2$	$eta_3$	$eta_4$	$\sigma_1^2$	$\sigma_2^2$	$\sigma_3^2$	$\sigma_4^2$
0.0	0.06	0.96	0.10	0.19	-0.01	-0.03	-0.02	0.02
au = 0.3	au = 0.3  0.06  0.26  0	0.19	0.19 0.13	0.09	0.09	0.03	0.06	
$\tau = 0.5$	0.02	0.20	0.10	0.07	0.01	-0.03	0.00	0.01
au = 0.5  0	0.05	0.20	0.10	0.07	0.11	0.08	0.07	0.07
$\tau = 0.7$	0.01	0.10	0.05	0.04	0.01	-0.05	0.01	0.02
au = 0.7	0.01 0.10	0.00	0.04	0.07	0.09	0.10	0.07	

Table 8: Mean Relative Squared Error (MRSE) values for fixed–effect estimates, and biases and standard errors for estimated variance components, obtained with 100 observations and 200 explanatory variables.

	$eta_1$	$eta_2$	$eta_3$	$eta_4$	$\sigma_1^2$	$\sigma_2^2$	$\sigma_3^2$	$\sigma_4^2$
- 0.2	$\tau = 0.3  0.06  0.15  0.18  0.$	0.10	-0.04	-0.05	0.01	-0.02		
$\gamma = 0.5$		0.18	0.18 0.10	0.04	0.09	0.05	0.05	
- 05	0.02	0.17	0.00	0.00 0.05	-0.05	0.00	-0.02	-0.01
$\gamma = 0.5$	= 0.5 0.03 0.17 0.09 (	0.05	0.06	0.19	0.04	0.04		
- 07	0.01	0.15	0.04	0.02	0.03	0.00	-0.01	-0.02
au = 0.7	0.01 0.15	0.04	0.03	0.08	0.14	0.05	0.05	

Some component planes are given on Figure 9 (150 explanatory variables) and Figure 10 (200 explanatory variables).



Figure 9: Component planes (1,2), (1,3) and (1,4) given by mixed-SCGLR for 100 observations, 150 explanatory variables and  $\tau = 0.3$ . The tuning parameter triplet (K, s, l) is set to (4, 0.5, 4).



Figure 10: Component planes (1, 2), (1, 3) and (1, 4) given by mixed-SCGLR for 100 observations, 200 explanatory variables, and  $\tau = 0.9$ . The tuning parameter triplet (K, s, l) is set to (4, 0.9, 4).

# References

- Anderson, D. A. and Aitkin, M. (1985). Variance Component Models with Binary Response: Interviewer Variability. Journal of the Royal Statistical Society, Series B (Methodological), 47(2):203–210.
- Bates, D., Mächler, M., Bolker, B., and Walker, S. (2015). Fitting linear mixed-effects models using lme4. *Journal of Statistical Software*, 67(1):1–48.
- Breslow, N. E. and Clayton, D. G. (1993). Approximate Inference in Generalized Linear Mixed Models. *Journal of the American Statistical Association*, 88(421):9–25.
- Bry, X., Trottier, C., Mortier, F., and Cornu, G. (2018). Component-based regularisation of a multivariate GLM with a thematic partitioning of the explanatory variables. *Statistical Modelling*. In press.
- Bry, X., Trottier, C., Mortier, F., Cornu, G., and Verron, T. (2014). Extending SCGLR to multiple regressor-groups: The Theme-SCGLR method. In *Proceedings of the eighth International Conference on Partial Least Squares and Related Methods*, Paris, France.
- Bry, X., Trottier, C., Mortier, F., Cornu, G., and Verron, T. (2016). The Multiple Facets of Partial Least Squares and Related Methods, chapter Supervised Component Generalized Linear Regression with Multiple Explanatory Blocks: THEME-SCGLR, pages 141–154. Springer International Publishing.
- Bry, X., Trottier, C., Verron, T., and Mortier, F. (2013). Supervised component generalized linear regression using a PLS-extension of the Fisher scoring algorithm. *Journal of Multivariate Analysis*, 119(4):47–60.
- Bry, X. and Verron, T. (2015). THEME: THEmatic Model Exploration through multiple co-structure maximization. *Journal of Chemometrics*, 29(12):637–647.
- Clayton, D. G. (1996). Generalized linear mixed models. In Markov chain Monte Carlo in practice, pages 275–301. Springer.

- Cornu, G., Mortier, F., Trottier, C., and Bry, X. (2018). SCGLR: Supervised Component Generalized Linear Regression. R package version 3.0.
- Eliot, M., Ferguson, J., Reilly, M. P., and Foulkes, A. S. (2011). Ridge Regression for Longitudinal Biomarker Data. *The International Journal of Biostatistics*, 7(1):1–11.
- Fahrmeir, L. and Tutz, G. (1994). Multivariate Statistical Modelling Based on Generalized Linear Models. Springer Series in Statistics. Springer-Verlag.
- Groll, A. (2017). glmmLasso: Variable Selection for Generalized Linear Mixed Models by L1-Penalized Estimation. R package version 1.5.1.
- Groll, A. and Tutz, G. (2014). Variable selection for generalized linear mixed models by  $L_1$ -penalized estimation. *Statistics and Computing*, 24(2):137–154.
- Hadfield, J. D. (2010). MCMC Methods for Multi-Response Generalized Linear Mixed Models: the MCMCglmm R Package. *Journal of Statistical Software*, 33(2):1–22.
- Henderson, C. R. (1975). Best Linear Unbiased Estimation and Prediction under a Selection Model. *Biometrics*, 31(2):423–447.
- Hoerl, A. E. and Kennard, R. W. (1970). Ridge Regression: Biased Estimation for Nonorthogonal Problems. *Technometrics*, 12(1):55–67.
- Knudson, C. (2016). Monte Carlo Likelihood Approximation for Generalized Linear Mixed Models. PhD thesis, University of Minnesota.
- McCullagh, P. and Nelder, J. (1989). Generalized Linear Models, Second Edition. Chapman & Hall/CRC Monographs on Statistics & Applied Probability. Taylor & Francis.
- McCulloch, C. E. (1997). Maximum Likelihood Algorithms for Generalized Linear Mixed Models. Journal of the American Statistical Association, 92(437):162–170.
- Nelder, J. A. and Wedderburn, R. W. M. (1972). Generalized Linear Models. Journal of the Royal Statistical Society. Series A (General), 135(3):370–384.

- Pinheiro, J. C. and Bates, D. M. (1995). Approximations to the Log-Likelihood Function in the Nonlinear Mixed-Effects Model. *Journal of Computational and Graphical Statistics*, 4(1):12–35.
- R Core Team (2017). R: A Language and Environment for Statistical Computing. R Foundation for Statistical Computing, Vienna, Austria.
- Schall, R. (1991). Estimation in Generalized Linear Models with Random Effects. Biometrika, 78(4):719–727.
- Schelldorfer, J., Meier, L., and Bühlmann, P. (2014). GLMMLasso: An Algorithm for High-Dimensional Generalized Linear Mixed Lodels Using l<sub>1</sub>-Penalization. Journal of Computational and Graphical Statistics, 23(2):460–477.
- Shun, Z. and McCullagh, P. (1995). Laplace Approximation of High Dimensional Integrals. Journal of the Royal Statistical Society, Series B (Methodological), 57(4):749–760.
- Tibshirani, R. (1996). Regression Shrinkage and Selection via the Lasso. Journal of the Royal Statistical Society, Series B (Methodological), 58(1):267–288.
- Zeger, S. L. and Karim, M. R. (1991). Generalized Linear Models With Random Effects; A Gibbs Sampling Approach. Journal of the American Statistical Association, 86(413):79– 86.
- Zhang, Y., Zhou, H., Zhou, J., and Sun, W. (2017). Regression Models for Multivariate Count Data. Journal of Computational and Graphical Statistics, 26(1):1–13.
- Zou, H. and Hastie, T. (2005). Regularization and variable selection via the elastic net. Journal of the Royal Statistical Society, Series B (Methodological), 67(2):301–320.