

# Introducing complex dependence structures into supervised component-based models

Jocelyn CHAUVET

work supervised by

**Catherine TROTTIER** and **Xavier BRY**

Montpellier

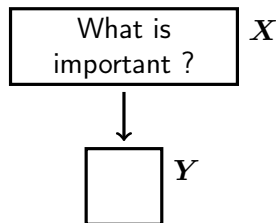
April 19<sup>th</sup>, 2019



↔ Why regularise a regression model ?

Fuzzy conceptual model, large amount of variables

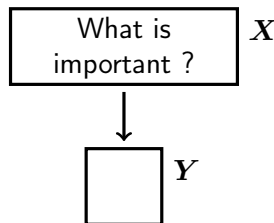
- ▶ Ill-conditioned matrix (almost singular)
  - ↔ Instability of coefficients
- ▶ High dimensional data ( $p > n$ )
  - ↔ Multicollinearity, singularity



## ↔ Why regularise a regression model ?

Fuzzy conceptual model, large amount of variables

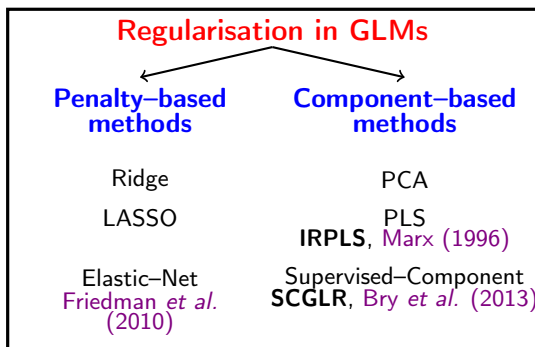
- ▶ Ill-conditioned matrix (almost singular)
  - ↔ Instability of coefficients
- ▶ High dimensional data ( $p > n$ )
  - ↔ Multicollinearity, singularity



### Regularisation (definition)

Introduction of additional criteria besides the Goodness-of-Fit in the estimation process in order to

- ▶ solve an ill-posed problem
- ▶ prevent overfitting



## Regularisation in GLMMs (GLMs + random effects)

Penalty-based  
methods:

LMM-Ridge  
Eliot *et al.* (2011)

GLMM-LASSO  
Groll and Tutz (2014)

### Regularisation in GLMs

Penalty-based  
methods

Component-based  
methods

Ridge

PCA

LASSO

PLS

IRPLS, Marx (1996)

Elastic-Net  
Friedman *et al.*  
(2010)

Supervised-Component  
SCGLR, Bry *et al.* (2013)

Component-based  
methods?

## Regularisation in GLMMs (GLMs + random effects)

Penalty-based  
methods:

LMM-Ridge  
Eliot *et al.* (2011)

GLMM-LASSO  
Groll and Tutz (2014)

### Regularisation in GLMs

Penalty-based  
methods

Ridge

LASSO

Elastic-Net  
Friedman *et al.*  
(2010)

Component-based  
methods

PCA

PLS

IRPLS, Marx (1996)

Supervised-Component  
SCGLR, Bry *et al.* (2013)

Component-based  
methods?

Extension of  
SCGLR  
to the GLMMs:  
Mixed-SCGLR

## Penalised log-likelihood

$$\ell(\beta; \mathbf{y}) - \lambda \text{pen}(\beta)$$

► **LASSO:**  $\text{pen}(\beta) = \|\beta\|_1$

- ↪ Sparse solutions
- ↪ Variable selection

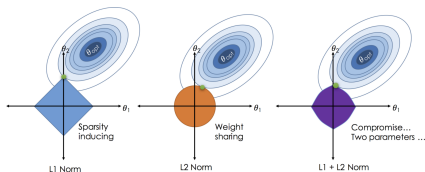
► **Ridge:**  $\text{pen}(\beta) = \|\beta\|_2^2$

- ↪ Shrinkage towards 0
- ↪ Reduce the estimator's variance

► **Elastic-net:**

$$\text{pen}(\beta) = (1 - \alpha)\|\beta\|_1 + \alpha\|\beta\|_2^2$$

- ↪ Variable selection
- ↪ Grouping effect



## Penalised log-likelihood

$$\ell(\beta; \mathbf{y}) - \lambda \text{pen}(\beta)$$

► **LASSO:**  $\text{pen}(\beta) = \|\beta\|_1$

- ↪ Sparse solutions
- ↪ Variable selection

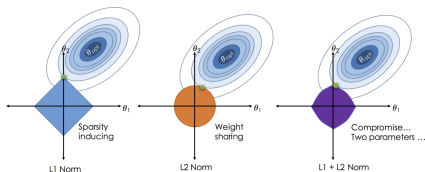
► **Ridge:**  $\text{pen}(\beta) = \|\beta\|_2^2$

- ↪ Shrinkage towards 0
- ↪ Reduce the estimator's variance

► **Elastic-net:**

$$\text{pen}(\beta) = (1 - \alpha)\|\beta\|_1 + \alpha\|\beta\|_2^2$$

- ↪ Variable selection
- ↪ Grouping effect



### Framework of interest

- Many highly correlated explanatory variables
- Proxies to latent phenomena to be found and interpreted



**Component-based approaches**



# Component-based regularisation

Components = synthetic variables  $\{f_h = X u_h \mid h = 1, \dots, H\}$

$$u_h = \begin{cases} \arg \max_{u \in \mathbb{R}^p} \text{crit}(u) \\ \text{w.r.t. } \|u\| = 1 \text{ and } Xu \perp Xu_1, \dots, Xu_{h-1} \end{cases}$$

# Component-based regularisation

Components = synthetic variables  $\{f_h = \mathbf{X}u_h \mid h = 1, \dots, H\}$

$$u_h = \begin{cases} \arg \max_{u \in \mathbb{R}^p} \text{crit}(u) \\ \text{w.r.t. } \|u\| = 1 \text{ and } \mathbf{X}u \perp \mathbf{X}u_1, \dots, \mathbf{X}u_{h-1} \end{cases}$$

## PCA

$$\text{crit}(u) = \underbrace{\|\mathbf{X}u\|_2^2}_{\text{Component Variance (CV)}}$$

↪ Information in  $\mathbf{X}$  ✓

↪ Link between  $\mathbf{X}$  and  $\mathbf{y}$  ✗

## PLS

$$\text{crit}(u) = \underbrace{\|\mathbf{X}u\|_2^2}_{\text{Comp. Var. (CV)}} \underbrace{\|\mathbf{y}\|_2^2 \cos^2(\mathbf{y}, \text{span}\{\mathbf{X}u\})}_{\text{Goodness-of-Fit (GoF)}}$$

↪ Information in  $\mathbf{X}$  ✓

↪ Link between  $\mathbf{X}$  and  $\mathbf{y}$  ✓

# Component-based regularisation

Components = synthetic variables  $\{f_h = X\mathbf{u}_h \mid h = 1, \dots, H\}$

$$\mathbf{u}_h = \begin{cases} \arg \max_{\mathbf{u} \in \mathbb{R}^p} \text{crit}(\mathbf{u}) \\ \text{w.r.t. } \|\mathbf{u}\| = 1 \text{ and } X\mathbf{u} \perp X\mathbf{u}_1, \dots, X\mathbf{u}_{h-1} \end{cases}$$

## PCA

$$\text{crit}(\mathbf{u}) = \underbrace{\|X\mathbf{u}\|_2^2}_{\text{Component Variance (CV)}}$$

↪ Information in  $X$  ✓

↪ Link between  $X$  and  $\mathbf{y}$  ✗

## PLS

$$\text{crit}(\mathbf{u}) = \underbrace{\|X\mathbf{u}\|_2^2}_{\text{Comp. Var. (CV)}} \underbrace{\|\mathbf{y}\|_2^2 \cos^2(\mathbf{y}, \text{span}\{X\mathbf{u}\})}_{\text{Goodness-of-Fit (GoF)}}$$

↪ Information in  $X$  ✓

↪ Link between  $X$  and  $\mathbf{y}$  ✓

## A "flexible PLS"

$$\text{crit}(\mathbf{u}) = [\text{CV}(\mathbf{u})]^s [\text{GoF}(\mathbf{u})]^{1-s}$$

$s \in [0, 1]$  a trade-off parameter

# Structural relevance

↪ Introduced by Bry and Verron (2015)

## Supervised Components via the Structural Relevance

$$\text{crit}(\mathbf{u}) = [\text{SR}(\mathbf{u})]^s [\text{GoF}(\mathbf{u})]^{1-s}$$

$$\text{SR}(\mathbf{u}) = \phi_l(\mathbf{u}) = \left( \sum_{j=1}^p [\text{cor}^2(\mathbf{X}\mathbf{u}, \mathbf{x}_j)]^l \right)^{\frac{1}{l}}$$

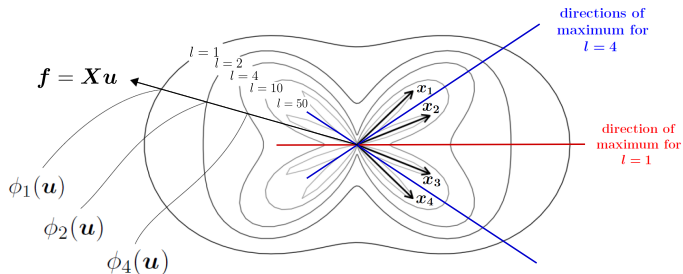
# Structural relevance

↪ Introduced by Bry and Verron (2015)

## Supervised Components via the Structural Relevance

$$\text{crit}(\mathbf{u}) = [\text{SR}(\mathbf{u})]^s [\text{GoF}(\mathbf{u})]^{1-s}$$

$$\text{SR}(\mathbf{u}) = \phi_l(\mathbf{u}) = \left( \sum_{j=1}^p [\text{cor}^2(\mathbf{X}\mathbf{u}, \mathbf{x}_j)]^l \right)^{\frac{1}{l}}$$

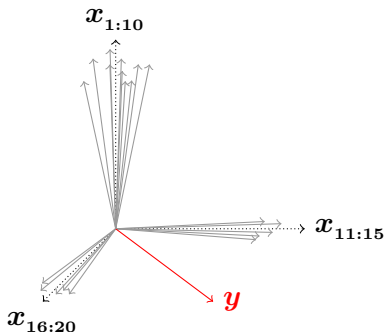


# Flash method-comparison

## A simple Gaussian model

▶  $y \sim \mathcal{N}_n(\mu = X\beta, \Sigma = I_n)$

▶  $X = \left[ \underbrace{x_1 \dots \dots \dots x_{10}}_{\substack{\text{large bundle} \\ \hookrightarrow \text{nuisance}}} \quad \underbrace{x_{11} \dots \dots \dots x_{15}}_{\substack{\text{small bundle} \\ \hookrightarrow \text{predicts } y}} \quad \underbrace{x_{16} \dots \dots \dots x_{20}}_{\substack{\text{small bundle} \\ \hookrightarrow \text{predicts } y}} \right]$



➡ PCR vs PLSR vs "Supervised Component Regression"

▶ **Non-independent observations**

↪ Grouped and panel data

⇒ From GLM to GLMM (use of random effects)

▶ **Multivariate framework**

↪ Several responses of various types  $\mathbf{Y} = [\mathbf{y}_1 \mid \dots \mid \mathbf{y}_q]$

▶ **Additional explanatory variables**

↪ with little redundancy

↪ Requiring no regularisation

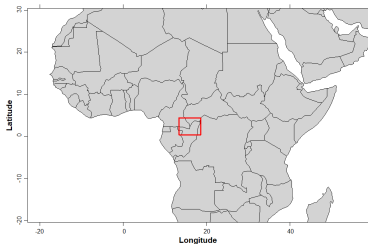
- 1 Supervised Component-based regularisation for multivariate GLMMs
  - The Congo-Basin floristic data
  - The mixed-SCGLR method
  - Simulation study
  - Results on the floristic data
- 2 Introducing a time-specific random effect
- 3 Perspectives



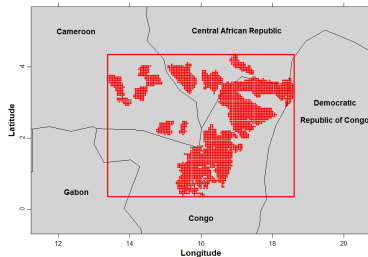
# Floristic data: a multivariate GLM ?

**Problem:** Model and predict the **abundance of tree species** in the tropical moist forest of the Congo–Basin

## The Congo–Basin



## 2615 land-plots



**Responses:**

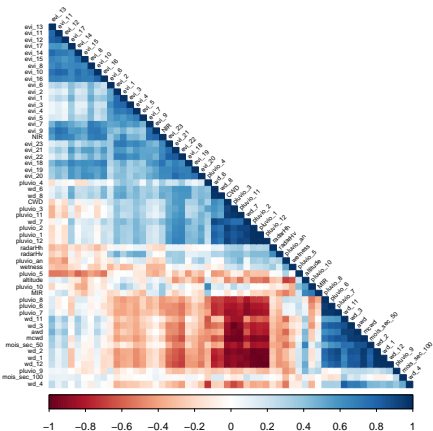
$q = 8$  abundances of selected tree species  
(i.e. **multivariate count** responses)



**Multivariate GLM**

# Regularisation needed

## Correlation-heatmap



## Explanatory variables:

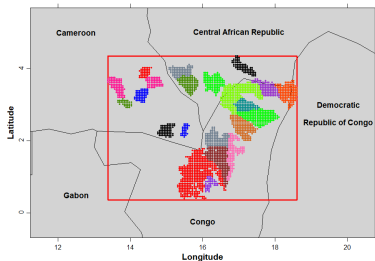
- ▶  $p = 56$  **highly correlated** environmental variables
  - ↳ Information **redundancy**
  - ↳ Model **instability**
- ▶  $r = 2$  **additional covariates** (geology and anthropogenic interference)



**Regularisation**  
via supervised components  
*common to all responses*

# Random effects needed: GLMM

## 22 forest concessions



### SCGLR:

- ▶ The land-plots are assumed **independent**...
- ▶ ...yet they are grouped into **22 concessions**

### Mixed-SCGLR:

- ▶ Within-group dependence modelled by a **random effect**



**Multivariate GLMM**

# Mixed-SCGLR (1)

## Notations:

- ▶  $\mathbf{Y}_{n \times q}$ : matrix of  $q$  response vectors  $\mathbf{y}_1, \dots, \mathbf{y}_q$
- ▶  $\mathbf{X}_{n \times p}$ : explanatory variables  $\mathbf{x}_1, \dots, \mathbf{x}_p$  (many, redundant)
- ▶  $\mathbf{A}_{n \times r}$ : additional covariates  $\mathbf{a}_1, \dots, \mathbf{a}_r$  (few, not redundant)
- ▶  $\mathbf{U}_{n \times N}$ : design matrix of the random effects

# Mixed-SCGLR (1)

## Notations:

- ▶  $\mathbf{Y}_{n \times q}$ : matrix of  $q$  response vectors  $\mathbf{y}_1, \dots, \mathbf{y}_q$
- ▶  $\mathbf{X}_{n \times p}$ : explanatory variables  $\mathbf{x}_1, \dots, \mathbf{x}_p$  (many, redundant)
- ▶  $\mathbf{A}_{n \times r}$ : additional covariates  $\mathbf{a}_1, \dots, \mathbf{a}_r$  (few, not redundant)
- ▶  $\mathbf{U}_{n \times N}$ : design matrix of the random effects

## Single component multivariate GLMM

For each  $k \in \{1, \dots, q\}$ ,

$$g_k(\mathbb{E}(\mathbf{Y}_k | \boldsymbol{\xi}_k)) = \boldsymbol{\eta}_k$$

$$\boldsymbol{\eta}_k = (\mathbf{X}\mathbf{u})\boldsymbol{\gamma}_k + \mathbf{A}\boldsymbol{\delta}_k + \mathbf{U}\boldsymbol{\xi}_k$$

$$\boldsymbol{\xi}_k \stackrel{\text{ind.}}{\sim} \mathcal{N}_N(\mathbf{0}, \sigma_k^2 \mathbf{I}_N), \text{ with } N \text{ the number of groups}$$

# Mixed-SCGLR (2)

Estimation method: Iterative procedure based on a linearisation of the model

↪ Pseudo-responses:  $z_k$

"Linearised" model

$$z_k = \underbrace{(X\mathbf{u})\gamma_k + A\delta_k + U\xi_k}_{\eta_k} + e_k \quad \text{with: } \begin{cases} \mathbb{E}(e_k | \xi_k) = 0 \\ \mathbb{V}(e_k | \xi_k) =: W_k^{-1} \end{cases}$$

# Mixed-SCGLR (2)

Estimation method: Iterative procedure based on a linearisation of the model

↪ Pseudo-responses:  $z_k$

## "Linearised" model

$$z_k = \underbrace{(X\mathbf{u})\gamma_k + A\delta_k + U\xi_k}_{\eta_k} + e_k \quad \text{with: } \begin{cases} \mathbb{E}(e_k | \xi_k) = 0 \\ \mathbb{V}(e_k | \xi_k) =: W_k^{-1} \end{cases}$$

## Alternate procedure

- (i) Given  $\gamma_k$ ,  $\delta_k$ ,  $\xi_k$  and  $\sigma_k^2$ , we compute the component  $f = X\mathbf{u}$
- (ii) Given  $\mathbf{u}$ , we estimate  $\gamma_k$ ,  $\delta_k$ ,  $\xi_k$  and  $\sigma_k^2$ 
  - ↪ Schall's algorithm, Henderson's system

# Mixed-SCGLR (3)

## Step (ii)

### Henderson's systems

↔ Given  $f = Xu$ , for each  $k \in \{1, \dots, q\}$ :

$$\begin{pmatrix} f^\top W_k f & f^\top W_k A & f^\top W_k U \\ A^\top W_k f & A^\top W_k A & A^\top W_k U \\ U^\top W_k f & U^\top W_k A & U^\top W_k U + D_k^{-1} \end{pmatrix} \begin{pmatrix} \gamma_k \\ \delta_k \\ \xi_k \end{pmatrix} = \begin{pmatrix} f^\top W_k z_k \\ A^\top W_k z_k \\ U^\top W_k z_k \end{pmatrix}$$

### Update variance components

$$\sigma_k^2 \leftarrow \frac{\xi_k^\top \xi_k}{N - \frac{1}{\sigma_k^2} \text{Trace} \left[ (U^\top W_k U + D_k^{-1})^{-1} \right]}$$



# Mixed-SCGLR (4)

## Step (i)

### Goodness-of-Fit

$$\psi(\mathbf{u}) = \sum_{k=1}^q \left\| \Pi_{\text{span}\{\mathbf{X}\mathbf{u}, \mathbf{A}\}} \mathbf{z}_k \right\|_{\mathbf{W}_k}^2$$

# Mixed-SCGLR (4)

## Step (i)

### Goodness-of-Fit

$$\psi(\mathbf{u}) = \sum_{k=1}^q \left\| \Pi_{\text{span}\{\mathbf{X}\mathbf{u}, \mathbf{A}\}} \mathbf{z}_k \right\|_{\mathbf{W}_k}^2$$

### Structural Relevance

$$\phi(\mathbf{u}) = \left( \sum_{j=1}^p \left[ \text{cor}^2(\mathbf{X}\mathbf{u}, \mathbf{x}_j) \right]^l \right)^{\frac{1}{l}}, \quad \text{with } l \in [1, +\infty)$$

# Mixed-SCGLR (4)

## Step (i)

### Goodness-of-Fit

$$\psi(\mathbf{u}) = \sum_{k=1}^q \left\| \Pi_{\text{span}\{\mathbf{X}\mathbf{u}, \mathbf{A}\}} \mathbf{z}_k \right\|_{\mathbf{W}_k}^2$$

### Structural Relevance

$$\phi(\mathbf{u}) = \left( \sum_{j=1}^p \left[ \text{cor}^2(\mathbf{X}\mathbf{u}, \mathbf{x}_j) \right]^l \right)^{\frac{1}{l}}, \quad \text{with } l \in [1, +\infty)$$

### Trade-off GoF/SR

$$\max \quad \left[ \psi(\mathbf{u}) \right]^{1-s} \left[ \phi(\mathbf{u}) \right]^s, \quad \text{with } s \in [0, 1]$$

$$\text{w.r.t. } \|\mathbf{u}\| = 1 \quad (\text{Identification constraint})$$

# Simulations (1)

Two random responses:  $Y = [y_1 | y_2]$

# Simulations (1)

Two random responses:  $Y = [y_1 \mid y_2]$

Fixed effects:

- ▶ 30 explanatory variables  $\mathcal{N}(0, 1)$ :

$$X = \left[ \underbrace{x_1 \dots \dots \dots x_{15}}_{\substack{\text{bundle } X_0 \text{ (large)} \\ \hookrightarrow \text{nuisance}}} \quad \underbrace{x_{16} \dots \dots \dots x_{25}}_{\substack{\text{bundle } X_1 \text{ (medium)} \\ \hookrightarrow \text{predicts } y_1}} \quad \underbrace{x_{26} \dots \dots \dots x_{30}}_{\substack{\text{bundle } X_2 \text{ (small)} \\ \hookrightarrow \text{predicts } y_2}} \right]$$

- ▶ Within each bundle:

$$\text{cor}(x_j, x_k) = \begin{cases} 1 & \text{if } j = k \\ \tau & \text{if } j \neq k \end{cases} \quad \text{with } \tau \in \{0.1, 0.3, 0.5, 0.7, 0.9\}$$

# Simulations (1)

Two random responses:  $Y = [y_1 \mid y_2]$

Fixed effects:

- ▶ 30 explanatory variables  $\mathcal{N}(0, 1)$ :

$$X = \left[ \underbrace{x_1 \dots x_{15}}_{\substack{\text{bundle } X_0 \text{ (large)} \\ \hookrightarrow \text{nuisance}}} \quad \underbrace{x_{16} \dots x_{25}}_{\substack{\text{bundle } X_1 \text{ (medium)} \\ \hookrightarrow \text{predicts } y_1}} \quad \underbrace{x_{26} \dots x_{30}}_{\substack{\text{bundle } X_2 \text{ (small)} \\ \hookrightarrow \text{predicts } y_2}} \right]$$

- ▶ Within each bundle:

$$\text{cor}(x_j, x_k) = \begin{cases} 1 & \text{if } j = k \\ \tau & \text{if } j \neq k \end{cases} \quad \text{with } \tau \in \{0.1, 0.3, 0.5, 0.7, 0.9\}$$

Random effects:

- ▶  $N = 10$  groups and  $R = 10$  units per group  
 $\hookrightarrow$  Design matrix  $U = I_N \otimes \mathbf{1}_R$

# Simulations (2)

## Model

$$\mathcal{M}: \begin{cases} y_1 = X\beta_1 + U\xi_1 + \varepsilon_1 \\ y_2 = X\beta_2 + U\xi_2 + \varepsilon_2 \end{cases} \quad \text{with}$$

$$\xi_k \sim \mathcal{N}_N(\mathbf{0}, \mathbf{I}_N) \quad \text{and} \quad \varepsilon_k \sim \mathcal{N}_{NR}(\mathbf{0}, \mathbf{I}_{NR})$$

100 samples of model  $\mathcal{M}$  for each value of  $\tau$

Comparison with:

### LMM-Ridge (2011)



Eliot et al.

↪ EM algorithm

↪ GCV at each step

### GLMM-LASSO (2014)



Groll, A. and Tutz, G.

↪ Laplace approximation

↪ Coordinate Gradient Descent

# Simulations (3)

$\tau$	LMM (No reg.)	GLMM-LASSO $\lambda_{\text{lasso}}^*$ (shrinkage)	LMM-Ridge $\lambda_{\text{ridge}}^*$ (shrinkage)	Mixed-SCGLR ( $l = 4$ ) $H^*$ (nb comp.)	$s^*$ (trade-off)
0.1		65	24	25	0.50
0.3		92	54	5	0.58
0.5		124	73	3	0.70
0.7		163	78	3	0.73
0.9		175	85	2	0.80

$\tau$	Ave $\left[ \max \left( \frac{\ \widehat{\beta}_1 - \beta_1\ _2^2}{\ \beta_1\ _2^2}, \frac{\ \widehat{\beta}_2 - \beta_2\ _2^2}{\ \beta_2\ _2^2} \right) \right]$				
0.1	0.12	0.05	0.08		0.12
0.3	0.33	0.12	0.13		0.10
0.5	0.61	0.20	0.16		0.07
0.7	1.32	0.25	0.20		0.06
0.9	4.62	0.26	0.31		0.05

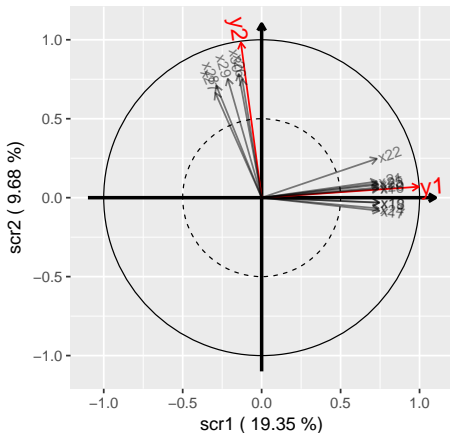


# Simulations (4)

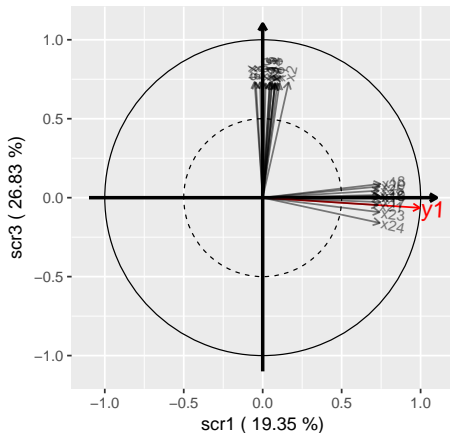
## Redundancy level

$$\tau = 0.5$$

### Component plane (1,2)



### Component plane (1,3)



# Simulations (5)

## Model

$$\mathcal{M}: \begin{cases} y_1 \sim \text{Ber}(p = \text{logit}^{-1} [X\beta_1 + U\xi_1]) \\ y_2 \sim \text{Poi}(\lambda = \exp [X\beta_2 + U\xi_2]), \end{cases}$$

Fixed-effect squared relative errors:

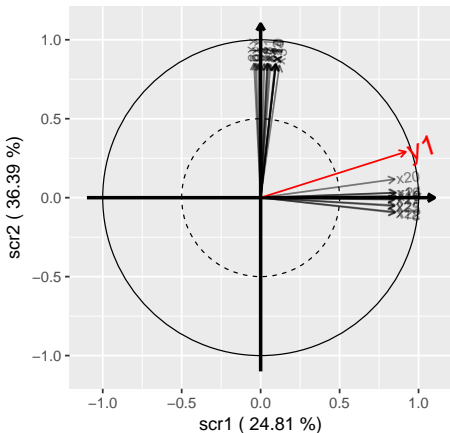
$\tau$	GLMM (no reg.)		GLMM-LASSO		mixed-SCGLR	
	Ber	Poi	Ber	Poi	Ber	Poi
0.1	316.48	0.54	8.61	0.30	14.71	0.46
0.3	398.78	0.64	9.23	0.36	7.21	0.21
0.5	576.68	0.87	14.48	0.44	2.01	0.09
0.7	886.04	1.28	17.37	0.47	1.50	0.07
0.9	2840.10	3.72	17.24	0.59	1.31	0.05

# Simulations (6)

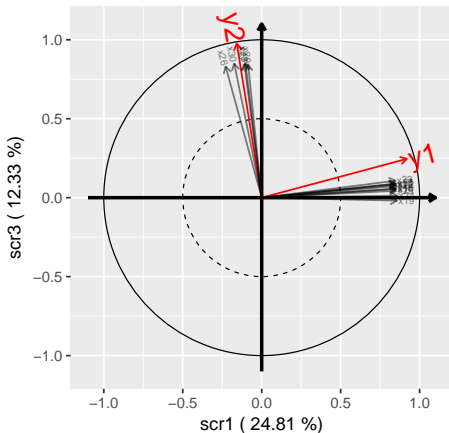
## Redundancy level

$$\tau = 0.9$$

Component plane (1,2)



Component plane (1,3)



# Results on the floristic data (1)

- ▶  $n = 2615$  land-plots, divided in  
 $N = 22$  forest concessions (considered as groups)
- ▶  $q = 8$  abundances of tree genera (responses  $\mathbf{Y}$ )
- ▶  $p = 56$  explanatory variables ( $\mathbf{X}$ )
- ▶  $r = 2$  additional covariates ( $\mathbf{A}$ )

# Results on the floristic data (1)

- ▶  $n = 2615$  land-plots, divided in  
 $N = 22$  forest concessions (considered as groups)
- ▶  $q = 8$  abundances of tree genera (responses  $\mathbf{Y}$ )
- ▶  $p = 56$  explanatory variables ( $\mathbf{X}$ )
- ▶  $r = 2$  additional covariates ( $\mathbf{A}$ )

## Model

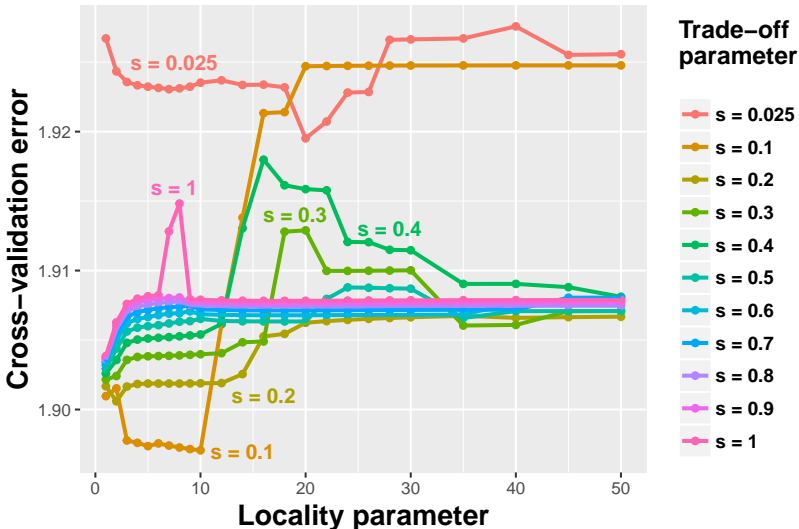
Abundance of tree species: count data

↔ Poisson regression with log link

$$\mathbf{y}_k \sim \mathcal{Poi}(\boldsymbol{\lambda} = \exp[\boldsymbol{\eta}_k])$$

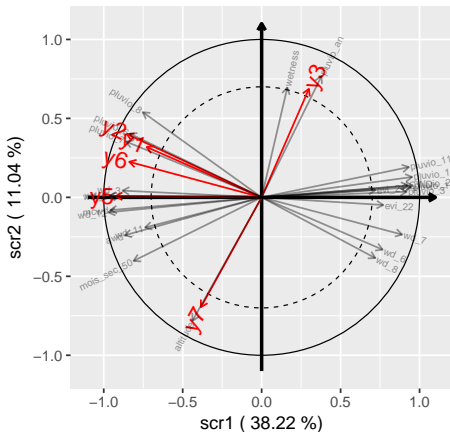
$$\boldsymbol{\eta}_k = \sum_{h=1}^H (\mathbf{X} \mathbf{u}_h) \gamma_{k,h} + \mathbf{A} \boldsymbol{\delta}_k + \mathbf{U} \boldsymbol{\xi}_k$$

# Results on the floristic data (2), $H^* = 4$

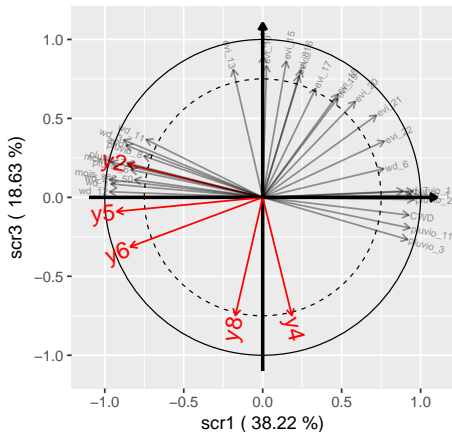


# Results on the floristic data (3)

Component plane (1,2)



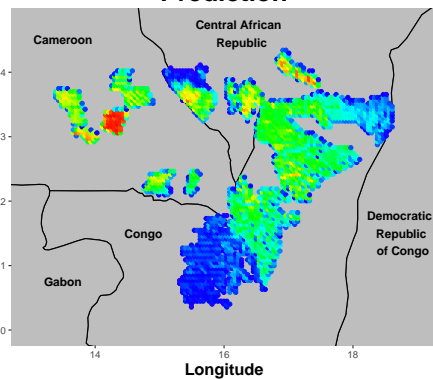
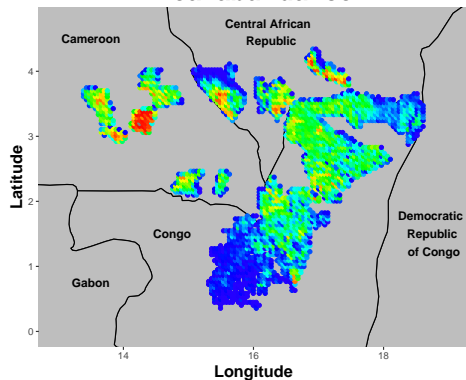
Component plane (1,3)



# Results on the floristic data (4)

## Real abundance

## Prediction





- 1 Supervised Component-based regularisation for multivariate GLMMs
- 2 Introducing a time-specific random effect
  - Model definition
  - Regularisation frameworks
    - Ridge-penalised EM
    - SC-regularised EM
  - Application to the GLMMs
  - Simulation study
- 3 Perspectives

# Framework of interest

Balanced panel data with:

- ▶  $N$  individuals...
- ▶ ...observed at the same  $R$  time-points

Notations:

- ▶  $\mathbf{y}_{NR \times 1}$ : response vector
- ▶  $\mathbf{X}_{NR \times p}$ : design matrix of the many and redundant explanatory variables

# Framework of interest

Balanced panel data with:

- ▶  $N$  individuals...
- ▶ ...observed at the same  $R$  time-points

Notations:

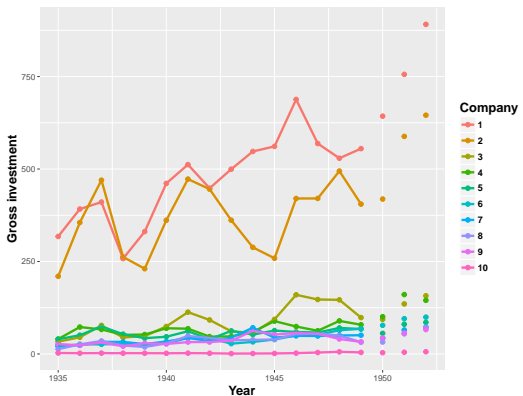
- ▶  $\mathbf{y}_{NR \times 1}$ : response vector
- ▶  $\mathbf{X}_{NR \times p}$ : design matrix of the many and redundant explanatory variables

## Difficulties

- ▶ High level of correlation among the explanatory variables  
⇒ Regularisation
- ▶ Individual- and time-specific effects  
⇒ Take into account the complex dependence structure

# The Grunfeld data

**Aim:** Predict the gross investments from 1950 onwards



Two elements to consider:

- ▶ The specific behaviour of each company
- ▶ The economic climate: latent phenomenon shared by all the companies which tends to persist over time

# Two-way random effects model

In general, we consider data with

- ▶ a **within-individual dependence**
  - ↪ **random effect with independent levels**
- ▶ a **time dependence**
  - ↪ **random effect with AR(1) levels**



**GLMM**

(in order to deal with non-Gaussian response)  
with both individual- and time-specific random effects

# The GLMM framework

$$\mathbf{y} = (y_{11}, y_{12}, \dots, y_{1R}, \\ y_{21}, y_{22}, \dots, y_{2R}, \dots \\ y_{N1}, y_{N2}, \dots, y_{NR})^T$$

$$g(\mathbb{E}(Y | \boldsymbol{\xi})) = \boldsymbol{\eta} = \mathbf{X}\boldsymbol{\beta} + \mathbf{U}_1\boldsymbol{\xi}_1 + \mathbf{U}_2\boldsymbol{\xi}_2$$

$$\mathbf{U}_1 = \mathbf{I}_N \otimes \mathbf{1}_R \quad \text{and} \quad \mathbf{U}_2 = \mathbf{1}_N \otimes \mathbf{I}_R$$

# The GLMM framework

$$\mathbf{y} = (y_{11}, y_{12}, \dots, y_{1R}, \\ y_{21}, y_{22}, \dots, y_{2R}, \dots \\ y_{N1}, y_{N2}, \dots, y_{NR})^\top$$

$$g(\mathbb{E}(Y | \boldsymbol{\xi})) = \boldsymbol{\eta} = \mathbf{X}\boldsymbol{\beta} + \mathbf{U}_1\boldsymbol{\xi}_1 + \mathbf{U}_2\boldsymbol{\xi}_2$$

$$\mathbf{U}_1 = \mathbf{I}_N \otimes \mathbf{1}_R \quad \text{and} \quad \mathbf{U}_2 = \mathbf{1}_N \otimes \mathbf{I}_R$$

- ▶ Individual-specific random effect:

$$\boldsymbol{\xi}_1 = (\xi_{11}, \xi_{12}, \dots, \xi_{1N})^\top \sim \mathcal{N}_N(\mathbf{0}, \sigma_1^2 \mathbf{A}_1), \quad \mathbf{A}_1 = \mathbf{I}_N$$

- ▶ Time-specific random effect:

$$\boldsymbol{\xi}_2 = (\xi_{21}, \xi_{22}, \dots, \xi_{2R})^\top \sim \mathcal{N}_R(\mathbf{0}, \sigma_2^2 \mathbf{A}_2(\rho))$$

$$\mathbf{A}_2(\rho) = \left( \frac{\rho^{|i-j|}}{1 - \rho^2} \right)_{1 \leq i, j \leq R}$$

- ▶  $\boldsymbol{\xi}_1 \perp \boldsymbol{\xi}_2$

# Estimation through penalised EM (1)

## Principle



**Green, P.J. (1990)** *On use of the EM for penalized likelihood estimation.*  
*Journal of the Royal Statistical Society. Series B (Methodological)*, 443-452.

$$\boldsymbol{\theta} = (\boldsymbol{\beta}, \sigma_1^2, \sigma_2^2, \rho), \boldsymbol{\xi} = (\boldsymbol{\xi}_1, \boldsymbol{\xi}_2)$$

$$\text{E} : \mathcal{Q}_{\text{pen}}(\boldsymbol{\theta} | \boldsymbol{\theta}^{[t]}) := \mathbb{E}_{\boldsymbol{\xi} | \mathbf{y}} \left[ \ell_{\text{pen}}^c(\boldsymbol{\theta}; \mathbf{y}, \boldsymbol{\xi}) | \boldsymbol{\theta}^{[t]} \right]$$

$$\text{M} : \boldsymbol{\theta}^{[t+1]} \leftarrow \arg \max_{\boldsymbol{\theta}} \mathcal{Q}_{\text{pen}}(\boldsymbol{\theta} | \boldsymbol{\theta}^{[t]})$$



# Estimation through penalised EM (1)

## Principle



**Green, P.J. (1990)** *On use of the EM for penalized likelihood estimation.*  
*Journal of the Royal Statistical Society. Series B (Methodological)*, 443-452.

$$\theta = (\beta, \sigma_1^2, \sigma_2^2, \rho), \quad \xi = (\xi_1, \xi_2)$$

$$E : Q_{\text{pen}}(\theta | \theta^{[t]}) := \mathbb{E}_{\xi | y} \left[ \ell_{\text{pen}}^c(\theta; y, \xi) | \theta^{[t]} \right]$$

$$M : \theta^{[t+1]} \leftarrow \arg \max_{\theta} Q_{\text{pen}}(\theta | \theta^{[t]})$$

## Usual penalised complete log-likelihood

$$\ell_{\text{pen}}^c(\theta; y, \xi) = \ell^c(\theta; y, \xi) - \lambda \text{pen}(\beta)$$

$$\text{pen}(\beta) = \begin{cases} \|\beta\|_1 \\ \|\beta\|_2^2 = \beta^T \beta \\ \alpha \|\beta\|_2^2 + (1 - \alpha) \|\beta\|_1, \quad 0 \leq \alpha \leq 1 \end{cases}$$

# Estimation through penalised EM (2)

## Ridge-based regularisation

$$\hookrightarrow \boldsymbol{\theta} = (\boldsymbol{\beta}, \sigma_1^2, \sigma_2^2, \rho), \boldsymbol{\xi} = (\boldsymbol{\xi}_1, \boldsymbol{\xi}_2)$$

$$\text{E} : Q_{\text{ridge}}(\boldsymbol{\theta}, \lambda \mid \boldsymbol{\theta}^{[t]}) := \mathbb{E}_{\boldsymbol{\xi} \mid \mathbf{y}} \left[ \ell^c(\boldsymbol{\theta}; \mathbf{y}, \boldsymbol{\xi}) - \lambda \boldsymbol{\beta}^T \boldsymbol{\beta} \mid \boldsymbol{\theta}^{[t]} \right]$$

$$\text{M} : \begin{cases} \lambda^{[t+1]} \leftarrow \text{GCV}^{[t+1]}(\lambda) \\ \boldsymbol{\theta}^{[t+1]} \leftarrow \arg \max_{\boldsymbol{\theta}} Q_{\text{ridge}}(\boldsymbol{\theta}, \lambda^{[t+1]} \mid \boldsymbol{\theta}^{[t]}) \end{cases}$$



Eliot, M., Ferguson, J., Reilly, M.P. and Foulkes, A.S. (2011) *Ridge Regression for Longitudinal Biomarker Data*. The International Journal of Biostatistics, **7**, 1–11.

# Estimation through SC-regularised EM

## Linear predictor

$$\eta = \begin{cases} (\mathbf{X}\mathbf{u})\gamma + \mathbf{U}_1\xi_1 + \mathbf{U}_2\xi_2 & \text{for a single component} \\ \sum_{h=1}^H (\mathbf{X}\mathbf{u}_h)\gamma_h + \mathbf{U}_1\xi_1 + \mathbf{U}_2\xi_2 & \text{for } H \text{ components} \end{cases}$$

# Estimation through SC-regularised EM

## Linear predictor

$$\eta = \begin{cases} (\mathbf{X}\mathbf{u})\gamma + \mathbf{U}_1\xi_1 + \mathbf{U}_2\xi_2 & \text{for a single component} \\ \sum_{h=1}^H (\mathbf{X}\mathbf{u}_h)\gamma_h + \mathbf{U}_1\xi_1 + \mathbf{U}_2\xi_2 & \text{for } H \text{ components} \end{cases}$$

## SC-based complete log-likelihood

$\theta = (\mathbf{u}, \gamma, \sigma_1^2, \sigma_2^2, \rho)$ , trade-off parameter  $s \in [0, 1]$

$$\ell_{\text{SC}}^c(\theta; \mathbf{y}, \xi) = (1-s) \ell^c(\theta; \mathbf{y}, \xi) + s \log[\phi(\mathbf{u})]$$

● **Complete log-likelihood** : measures the fit of the model (based on component  $\mathbf{f} = \mathbf{X}\mathbf{u}$ ) to the data

● **Structural relevance criterion** : measures the closeness of component  $\mathbf{f}$  to the strongest structures of  $\mathbf{X}$

# Ridge– versus SC–EM

## Ridge–based penalisation

$$\ell^c(\boldsymbol{\theta}; \mathbf{y}, \boldsymbol{\xi}) - \lambda \boldsymbol{\beta}^\top \boldsymbol{\beta}$$

- ▶ **Penalises** the "large" coefficients
- ▶ Sees the high correlations among the explanatory variables as **pure nuisance**
- ▶  $\eta$  **hard to interpret**

## Component–based regularisation

$$\ell^c(\boldsymbol{\theta}; \mathbf{y}, \boldsymbol{\xi}) + \frac{s}{1-s} \log [\phi(\mathbf{u})]$$

- ▶ **Gives a bonus** to the most interpretable bundles in  $\mathbf{X}$
- ▶ **Takes advantage** of the high correlations among the explanatory variables
- ▶  $\eta$  **easier to interpret** through decomposition on components

# Application to the GLMMs

## LINEARISATION step

- ▶ Given  $\boldsymbol{\mu}_i := \mathbb{E}(\mathbf{Y}_i | \boldsymbol{\xi})$ , the working variable  $\mathbf{z}_i$  writes

$$\mathbf{z}_i = \underbrace{g(\boldsymbol{\mu}_i)}_{\boldsymbol{\eta}_i} + (\mathbf{y}_i - \boldsymbol{\mu}_i)g'(\boldsymbol{\mu}_i)$$

- ▶ Linearised model:

$$\mathcal{M}: \mathbf{z} = \underbrace{\mathbf{X}\boldsymbol{\beta} + \mathbf{U}_1\boldsymbol{\xi}_1 + \mathbf{U}_2\boldsymbol{\xi}_2}_{\boldsymbol{\eta}} + \mathbf{e}, \quad \text{with } \mathbb{V}(\mathbf{e}) = \mathbf{W}$$

## ESTIMATION step

- ▶ Instead of the classical Henderson's systems, we propose a Penalised/Regularised EM algorithm on  $\mathcal{M}$

# Ridge-EM for GLMM-AR<sub>1</sub>

$$\theta = (\beta, \sigma_1^2, \sigma_2^2, \rho)$$

## Linearised model

$$\mathcal{M}^{[t]}: z^{[t]} = X\beta + U_1\xi_1 + U_2\xi_2 + e^{[t]}, \quad \text{with } \mathbb{V}(e^{[t]}) = W^{[t]}$$

## Ridge estimation

$$\mathbf{E}: Q_{\text{ridge}}(\theta, \lambda | \theta^{[t]}) := \mathbb{E}_{\xi | z^{[t]}} \left[ \ell^c(\theta; z^{[t]}, \xi) - \lambda \beta^T \beta | \theta^{[t]} \right]$$

$$\mathbf{M}: \begin{cases} \lambda^{[t+1]} \leftarrow \text{GCV}^{[t+1]}(\lambda) \\ \theta^{[t+1]} \leftarrow \arg \max_{\theta} Q_{\text{ridge}}(\theta, \lambda^{[t+1]} | \theta^{[t]}) \end{cases}$$

## Update

Compute  $\xi^{[t+1]}, z^{[t+1]}, W^{[t+1]}$

# SC-EM for GLMM-AR<sub>1</sub>

$$\boldsymbol{\theta} = (\mathbf{u}, \gamma, \sigma_1^2, \sigma_2^2, \rho)$$

## Linearised model

$$\mathcal{M}^{[t]}: \mathbf{z}^{[t]} = (\mathbf{X}\mathbf{u})\gamma + \mathbf{U}_1\xi_1 + \mathbf{U}_2\xi_2 + \mathbf{e}^{[t]}, \quad \text{with } \mathbb{V}(\mathbf{e}^{[t]}) = \mathbf{W}^{[t]}$$

## SC-estimation

$$\mathbf{E}: \mathcal{Q}_{\text{SC}}(\boldsymbol{\theta} | \boldsymbol{\theta}^{[t]}) := \mathbb{E}_{\boldsymbol{\xi} | \mathbf{z}^{[t]}} \left[ (1-s)\ell^c(\boldsymbol{\theta}; \mathbf{z}^{[t]}, \boldsymbol{\xi}) + s \log[\phi(\mathbf{u})] | \boldsymbol{\theta}^{[t]} \right]$$

$$\mathbf{M}: \begin{cases} \sigma_1^{2[t+1]}, \sigma_2^{2[t+1]}, \rho^{[t+1]} \text{ computed as previously} \\ \mathbf{u}^{[t+1]} \leftarrow \arg \max_{\|\mathbf{u}\|=1} \tilde{\mathcal{Q}}_{\text{SC}}(\mathbf{u}, \gamma^{[t]} | \boldsymbol{\theta}^{[t]}) \\ \gamma^{[t+1]} \leftarrow \arg \max_{\gamma} \tilde{\mathcal{Q}}_{\text{SC}}(\mathbf{u}^{[t+1]}, \gamma | \boldsymbol{\theta}^{[t]}) \end{cases}$$

## Update

$$\text{Compute } \boldsymbol{\xi}^{[t+1]}, \mathbf{z}^{[t+1]}, \mathbf{W}^{[t+1]}$$



# Simulations (1)

## Poisson regression with log link

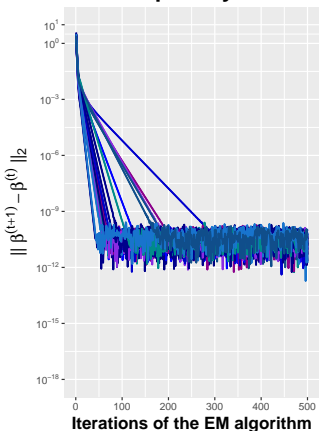
$$\blacktriangleright y \sim \text{Poi}(\lambda = \exp [X\beta + U_1\xi_1 + U_2\xi_2])$$

$$\blacktriangleright X = \left[ \underbrace{x_1 \dots \dots \dots x_{10}}_{\substack{\text{large bundle} \\ \hookrightarrow \text{noise}}} \quad \underbrace{x_{11} \dots \dots \dots x_{15}}_{\substack{\text{small bundle} \\ \hookrightarrow \text{predicts } y}} \quad \underbrace{x_{16} \dots \dots \dots x_{20}}_{\substack{\text{small bundle} \\ \hookrightarrow \text{predicts } y}} \right]$$

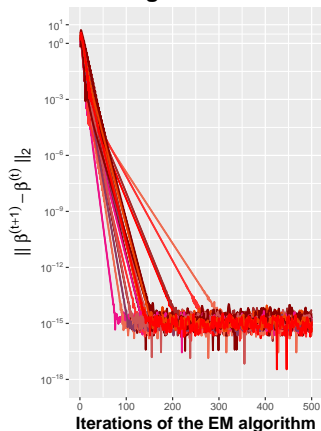
# Simulations (2)

How does convergence go?

**Ridge  
penalty**



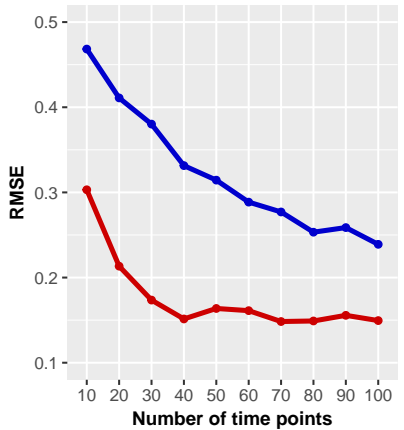
**Supervised Component  
regularisation**



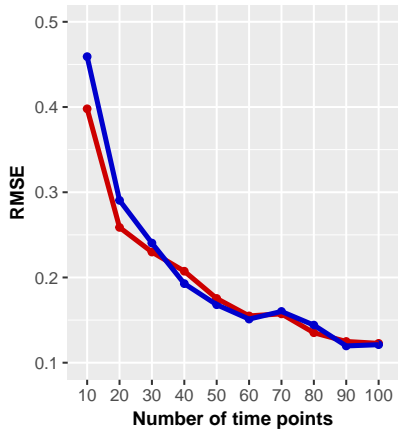
# Simulations (3)

## Accuracy of the estimates

Fixed effects parameter  $\beta$



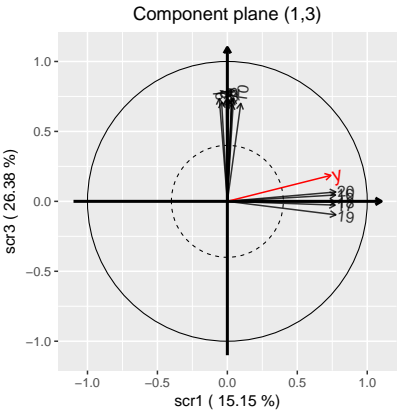
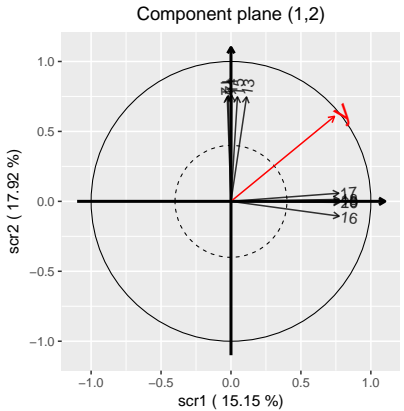
Autocorrelation  $\rho$



—●— SC —●— ridge

# Simulations (4)

## Power for model interpretation



## Powerful trade-off between

- ▶ multivariate GLMM
- ▶ component-based regularisation

## Model interpretation ↗

- ▶ Mixed-SCGLR provides graphical diagnoses (component planes)
- ▶ reveals the multidimensional explanatory and predictive structures

## Estimate-accuracy ↗

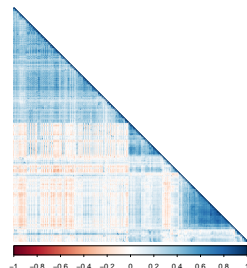
## Mixed-SCGLR now suitable for

- ▶ grouped data
- ▶ panel data

## Major depressive disorders and grey-matter volume reduction

- ▶ 15 years follow-up neuropsychiatric study including **636 participants**
- ▶ **Repeated binary response:** Depressive or not
- ▶ **528 explanatory variables:** thickness, area, volume and curvature of several brain areas
- ▶ **Additional covariates:** Gender, age. . .

Correlation-heatmap



## Sparse Supervised Component $\{f_h = X u_h \mid h = 1, \dots, K\}$

- ▶ Variable selection via supervised components
  - ↪ Supervised components that are weighted combination of only a few explanatory variables
- ▶ Idea: add a **sparsity constraint** on the input variables
  - ↪ Generic program:

$$u_h = \begin{cases} \arg \max_{u \in \mathbb{R}^p} s \log [\phi(u)] + (1 - s) \log [\psi(u)] - \lambda \|u\|_1 \\ \|u\| = 1 \text{ and } Xu \perp Xu_1, \dots, Xu_{h-1} \end{cases}$$

## Spatial correlation modelling

$$\xi \sim \mathcal{N}_q(\mathbf{0}, \varsigma^2 \mathbf{B}^{-1}), \text{ where}$$

- ▶  $\varsigma^2$  is the unknown spatial-specific variance component, and
- ▶  $\mathbf{B} = \mathbf{I} - \rho \mathbf{A}$ , where  $\rho$  is the unknown autoregressive spatial parameter and  $\mathbf{A} = (a_{ij})_{1 \leq i, j \leq q}$  is the adjacency matrix defined by

$$a_{ij} = \begin{cases} 1 & \text{if } j \text{ is adjacent to } i \\ 0 & \text{otherwise.} \end{cases}$$

## A general covariance structure on $\mathbf{Y} = [\mathbf{y}_1 \mid \dots \mid \mathbf{y}_q]$ ?

- ▶ e.g. modelling species interactions and competition
  - ↪ Congo-Basin floristic data: some species tend to appear together
  - ↪ Other species are antagonistic



# References



**Bry, X., Verron, T. (2015).** *THEME: THEmatic model exploration through multiple co-structure maximization.* Journal of Chemometrics, **29**, 637-647.



**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**, 47-60.

+ Package R : **SCGLR**

<https://github.com/SCnext/SCGLR>



**Chauvet, J., Bry, X., Trottier, C. (2019).** *Component-based regularisation of multivariate generalised mixed models.* Journal of Computational and Graphical Statistics, *in press*.

+ Package R : **mixedSCGLR**

<https://github.com/SCnext/mixedSCGLR>



**Chauvet, J., Bry, X., Trottier, C. (2019).** *Regularisation of GLMMs with an autoregressive random time-specific effect.* *in progress*.



**Eliot, M., Ferguson, J., Reilly, M.P. and Foulkes, A.S. (2011)** *Ridge Regression for Longitudinal Biomarker Data.* The International Journal of Biostatistics, **7**, 1-11.



**Green, P.J. (1990)** *On use of the EM for penalized likelihood estimation.* Journal of the Royal Statistical Society. Series B (Methodological), 443-452.



**Marx, B. D. (1996)** *Iteratively reweighted partial least squares estimation for generalized linear regression.* Technometrics, **38**, 4, 374-381.

# High dimensional data (1)

- ▶ 100 observations
- ▶ 150 explanatory variables:  $X = [X_0 \mid X_1 \mid X_2 \mid X_3]$ 
  - ↪ 60 variables in  $X_0$ , 45 variables in  $X_1$ ,  
30 variables in  $X_2$ , 15 variables in  $X_3$

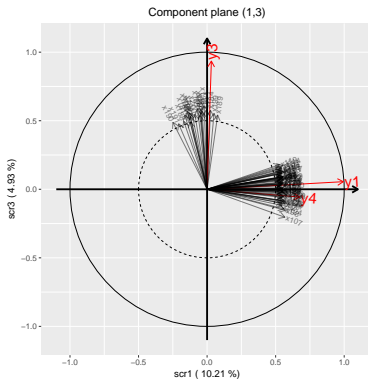
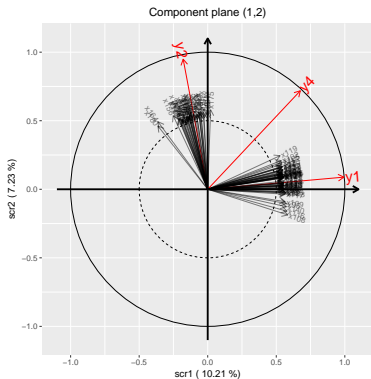
- ▶ Responses

$$\begin{cases} y_1 \sim \mathcal{N}_n(\mu = X\beta_1 + U\xi_1, \Sigma = I_n) \\ y_2 \sim \mathcal{B}(p = \text{logit}^{-1}[X\beta_2 + U\xi_2]) \\ y_3 \sim \text{Bin}(\text{trials} = 30 \mathbf{1}_n, p = \text{logit}^{-1}[X\beta_3 + U\xi_3]) \\ y_4 \sim \mathcal{P}(\lambda = \exp[X\beta_4 + U\xi_4]). \end{cases}$$

- ▶  $X_0$  = nuisance bundle  
 $y_1$  predicted only by  $X_1$ ,  $y_2$  only by  $X_2$ ,  $y_3$  only by  $X_3$ ,  
 $y_4$  by both  $X_2$  and  $X_3$ .

# High dimensional data (2)

## Results



# High dimensional data (3)

First idea: replace  $X$  with the matrix  $C$  of its principal components associated with non-negligible eigenvalues

- ▶  $C = XV$ , with  $V$  the matrix of unit-eigenvectors
- ▶ Modified GoF and SR criteria

$$\tilde{\psi}(\mathbf{u}) = \sum_{k=1}^q \left\| \Pi_{\text{span}\{\mathbf{C}\mathbf{u}, \mathbf{A}\}} \mathbf{z}_k \right\|_{\mathbf{W}_k}^2$$
$$\tilde{\phi}(\mathbf{u}) = \left( \sum_{j=1}^p \left[ \text{cor}^2(\mathbf{C}\mathbf{u}, \mathbf{x}_j) \right]^l \right)^{\frac{1}{l}}$$

- ▶ Maximisation program

$$\begin{cases} \max & s \log [\tilde{\phi}(\mathbf{u})] + (1 - s) \log [\tilde{\psi}(\mathbf{u})] \\ \text{subject to} & \mathbf{u}^\top \mathbf{C}^\top \mathbf{P} \mathbf{C} \mathbf{u} = 1, \mathbf{P} = n^{-1} \mathbf{I}_n \end{cases}$$

# High dimensional data (4)

**Better idea:**

- ▶ Preserve the GoF and SR criteria

$$\psi(\mathbf{u}) = \sum_{k=1}^q \left\| \Pi_{\text{span}\{\mathbf{X}\mathbf{u}, \mathbf{A}\}} \mathbf{z}_k \right\|_{\mathbf{W}_k}^2$$
$$\phi(\mathbf{u}) = \left( \sum_{j=1}^p \left[ \text{cor}^2(\mathbf{X}\mathbf{u}, \mathbf{x}_j) \right]^l \right)^{\frac{1}{l}}$$

- ▶ Modify the norm-constraint

$$\begin{cases} \max & s \log[\phi(\mathbf{u})] + (1-s) \log[\psi(\mathbf{u})] \\ \text{subject to} & \mathbf{u}^\top [\tau \mathbf{I} + (1-\tau) \mathbf{X}^\top \mathbf{P} \mathbf{X}] \mathbf{u} = 1, \mathbf{P} = n^{-1} \mathbf{I}_n \end{cases}$$

# The PING algorithm

- ▶ Generic program:

$$\begin{cases} \max & \mathcal{J}(\mathbf{v}) \\ \text{subject to} & \mathbf{v}^\top \mathbf{v} = 1 \text{ and } \mathbf{\Delta}^\top \mathbf{v} = \mathbf{0}. \end{cases}$$

- ▶ Direction of ascent:

$$\mathbf{v}^{[t+1]} = \frac{\Pi_{\text{span}\{\mathbf{\Delta}\}^\perp} \Gamma(\mathbf{v}^{[t]})}{\left\| \Pi_{\text{span}\{\mathbf{\Delta}\}^\perp} \Gamma(\mathbf{v}^{[t]}) \right\|}, \text{ with } \Gamma(\mathbf{v}) = \nabla \mathcal{J}(\mathbf{v})$$

# The PING algorithm

## The Projected Iterated Normed Gradient Algorithm

**while** *convergence of  $v$  non reached* **do**

$$m \leftarrow \frac{\Pi_{\text{span}\{\Delta\}^\perp} \Gamma(v^{[t]})}{\left\| \Pi_{\text{span}\{\Delta\}^\perp} \Gamma(v^{[t]}) \right\|}$$

**while**  $\mathcal{J}(m) < \mathcal{J}(v^{[t]})$  **do**

$$m \leftarrow \frac{v^{[t]} + m}{\|v^{[t]} + m\|}$$

**end**

$$v^{[t+1]} \leftarrow m$$

$$t \leftarrow t + 1$$

**end**