

Algorithme EM régularisé pour données longitudinales

Jocelyn CHAUVET
Catherine TROTTIER, Xavier BRY

RJS 2017, Porquerolles
6 avril 2017



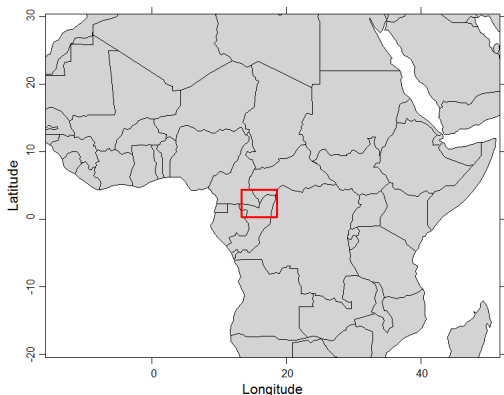
- 1 Régression Linéaire Généralisée sur Composantes Supervisées (SCGLR)
- 2 Ridge GL2M pour données longitudinales
- 3 Bilan et perspectives

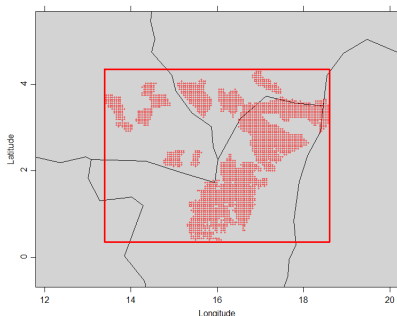
1 Régression Linéaire Généralisée sur Composantes Supervisées (SCGLR)

- Motivation
- Modèle et méthode d'estimation (première composante)
- Diagnostics graphiques sur données réelles

Problème

Modéliser et prédire les **distributions d'abondance d'espèces d'arbres** dans les forêts tropicales du bassin du Congo





Sur chaque parcelle :

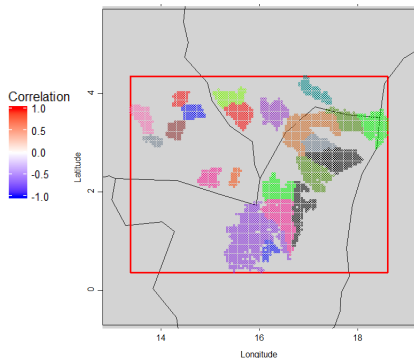
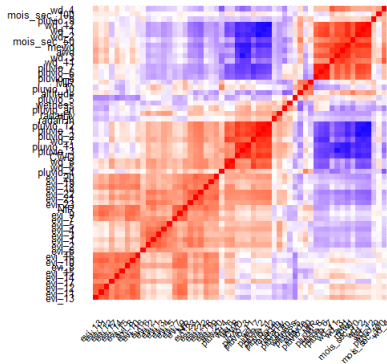
- ▶ $q = 94$ espèces d'arbres à prédire (abondance : données de comptage)

Dans le but de modéliser ces distributions d'espèces :

- ▶ $p = 56$ variables explicatives

Difficultés

- Variables explicatives très corrélées
↳ **Régularisation** nécessaire
- Observations géo-référencées
↳ **Structure de dépendance** plus complexe



SCGLR

- Gère les fortes redondances au sein des variables explicatives
 - ↪ **Régularisation par construction de composantes supervisées**

Mixed-SCGLR

- Prends en compte la structure de dépendance
 - ↪ Dépendance intra-groupe modélisée par un **effet aléatoire**
 - ↪ **GLMM multivarié**

- 1 Régression Linéaire Généralisée sur Composantes Supervisées (SCGLR)
 - Motivation
 - Modèle et méthode d'estimation (première composante)
 - Diagnostics graphiques sur données réelles

Notations :

- ▶ $\mathbf{Y}_{n \times q}$: matrice des q réponses y^1, \dots, y^q
- ▶ $\mathbf{X}_{n \times p}$: variables explicatives (redondantes)
- ▶ $\mathbf{U}_{n \times N}$: matrice de design des effets aléatoires

Notations :

- ▶ $Y_{n \times q}$: matrice des q réponses y^1, \dots, y^q
- ▶ $X_{n \times p}$: variables explicatives (redundantes)
- ▶ $U_{n \times N}$: matrice de design des effets aléatoires

Définition des prédicteurs linéaires : $\forall k \in \{1, \dots, q\}$,

$$\eta^k = g \left(\mathbb{E} \left(y^k \mid \xi_k \right) \right)$$

Modélisation classique

$$\eta^k = X\beta_k + U\xi_k$$

Notations :

- ▶ $Y_{n \times q}$: matrice des q réponses y^1, \dots, y^q
- ▶ $X_{n \times p}$: variables explicatives (redundantes)
- ▶ $U_{n \times N}$: matrice de design des effets aléatoires

Définition des prédicteurs linéaires : $\forall k \in \{1, \dots, q\}$,

$$\eta^k = g \left(\mathbb{E} \left(y^k \mid \xi_k \right) \right)$$

Régularisation

$$\eta^k = (X\mathbf{u})\gamma_k + U\xi_k$$

Comment calcule-t-on la composante $f = Xu$?
Bry et Verron (2015)

$$\max_{\|u\|=1} \left\{ C(u) = [\psi(u)]^{1-s} \times [\phi(u)]^s \right\}$$

Comment calcule-t-on la composante $f = Xu$?

Bry et Verron (2015)

$$\max_{\|u\|=1} \left\{ C(u) = [\psi(u)]^{1-s} \times [\phi(u)]^s \right\}$$

- Critère de qualité d'ajustement : mesure le degré d'ajustement de l'ensemble des prédicteurs linéaires η^k sur la composante $f = Xu$

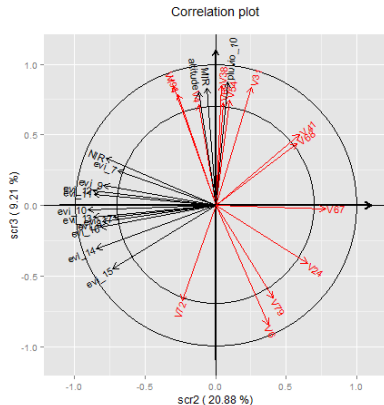
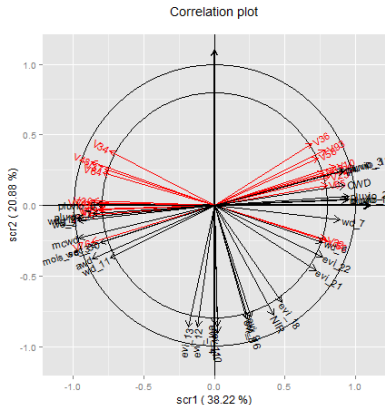
Comment calcule-t-on la composante $f = Xu$? Bry et Verron (2015)

$$\max_{\|u\|=1} \left\{ C(u) = [\psi(u)]^{1-s} \times [\phi(u)]^s \right\}$$

- Critère de qualité d'ajustement : mesure le degré d'ajustement de l'ensemble des prédicteurs linéaires η^k sur la composante $f = Xu$
- Critère de pertinence structurelle : mesure à quel point la composante est proche des structures fortes de X

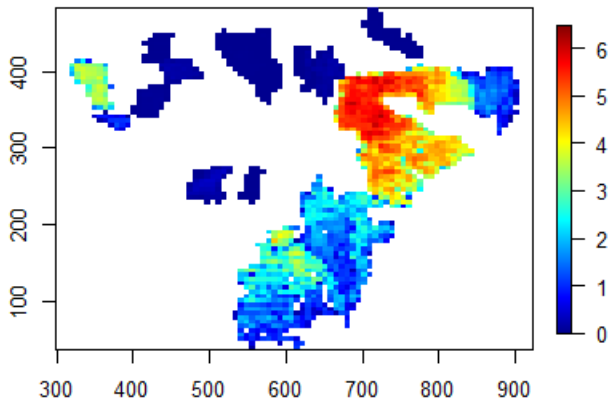
- 1 Régression Linéaire Généralisée sur Composantes Supervisées (SCGLR)
 - Motivation
 - Modèle et méthode d'estimation (première composante)
 - Diagnostics graphiques sur données réelles

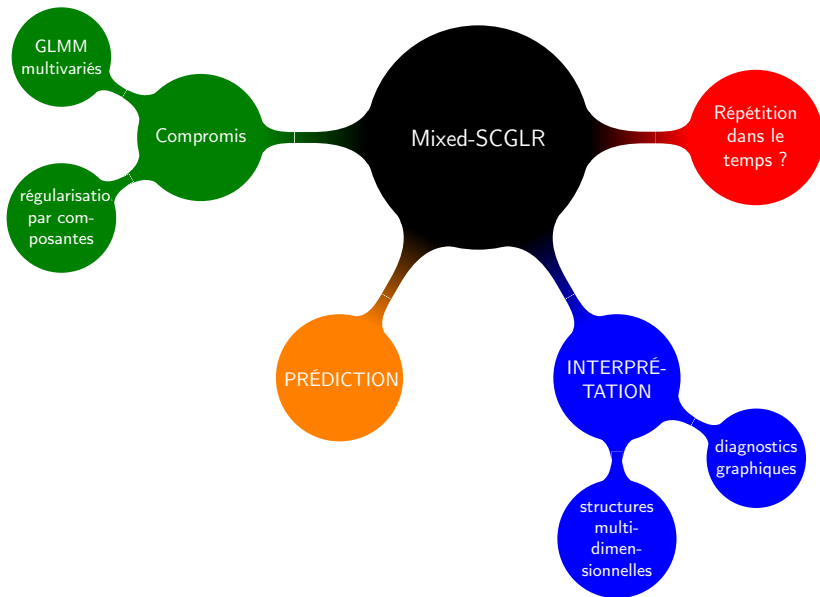
Qualités interprétatives Plans factoriels (1,2) & (2,3)



Cartes prédictives

Prédictions Espèce 17





2 Ridge GL2M pour données longitudinales

- L'existant : L2M-Ridge
- Adaptation aux données longitudinales : ajout d'un effet aléatoire temporel
- Extension aux GL2M
- Simulations et données réelles

L2M pour données répétées : $y = X\beta + U\xi + \varepsilon$

- ▶ $y = (y_{11}, \dots, y_{1R}, y_{21}, \dots, y_{2R}, \dots, y_{N1}, \dots, y_{NR})^T$
- ▶ β vecteur des effets fixes et X leur matrice de design
- ▶ $\xi = (\xi_1, \xi_2, \dots, \xi_N)^T$ vecteur de l'effet aléatoire "individu" et $U = I_N \otimes \mathbf{1}_R$ matrice de design correspondante
- ▶ $\varepsilon \sim \mathcal{N}(0, \sigma_0^2 I_{NR}) \quad \perp \quad \xi \sim \mathcal{N}(0, \sigma_1^2 I_N)$

L2M pour données répétées : $y = X\beta + U\xi + \varepsilon$

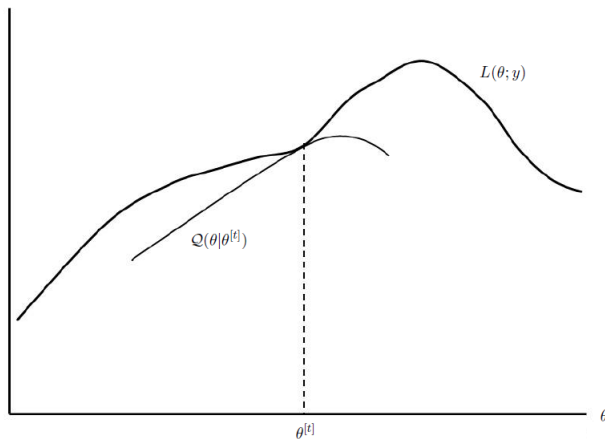
- ▶ $y = (y_{11}, \dots, y_{1R}, y_{21}, \dots, y_{2R}, \dots, y_{N1}, \dots, y_{NR})^T$
- ▶ β vecteur des effets fixes et X leur matrice de design
- ▶ $\xi = (\xi_1, \xi_2, \dots, \xi_N)^T$ vecteur de l'effet aléatoire "individu" et $U = I_N \otimes \mathbf{1}_R$ matrice de design correspondante
- ▶ $\varepsilon \sim \mathcal{N}(0, \sigma_0^2 I_{NR}) \perp \xi \sim \mathcal{N}(0, \sigma_1^2 I_N)$

Estimation classique par algorithme EM, $\theta = (\beta^T, \sigma_0^2, \sigma_1^2)^T$

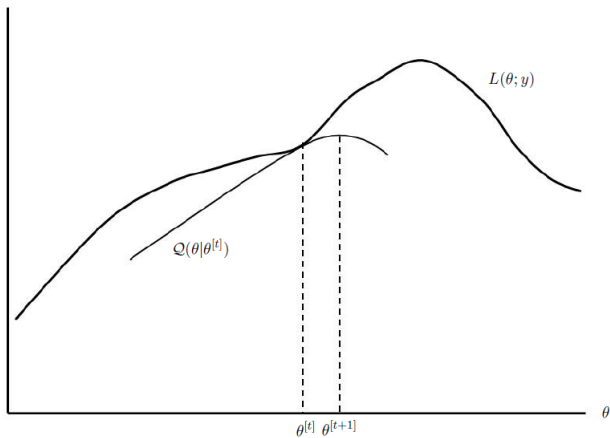
$$\mathbf{E} : Q(\theta | \theta^{[t]}) := \mathbb{E}_{\xi|y} [L(\theta; y, \xi) | \theta^{[t]}]$$

$$\mathbf{M} : \theta^{[t+1]} \leftarrow \arg \max_{\theta} Q(\theta | \theta^{[t]})$$

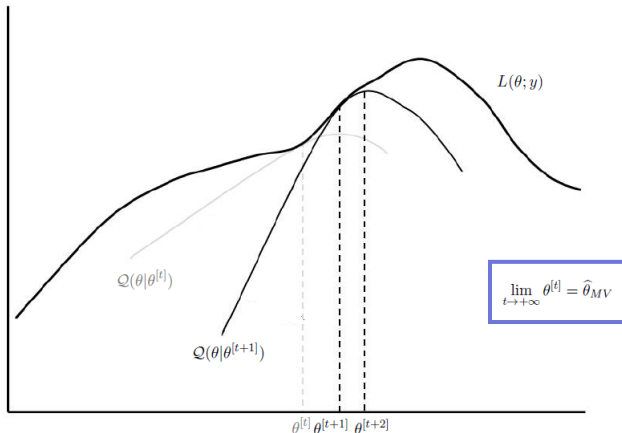
Algorithme EM classique



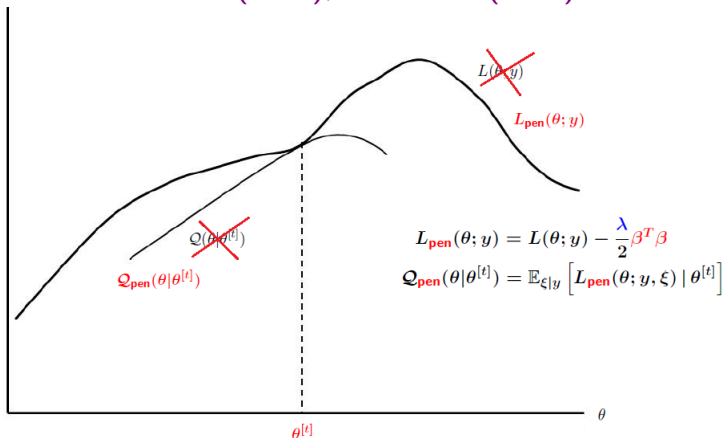
Algorithme EM classique



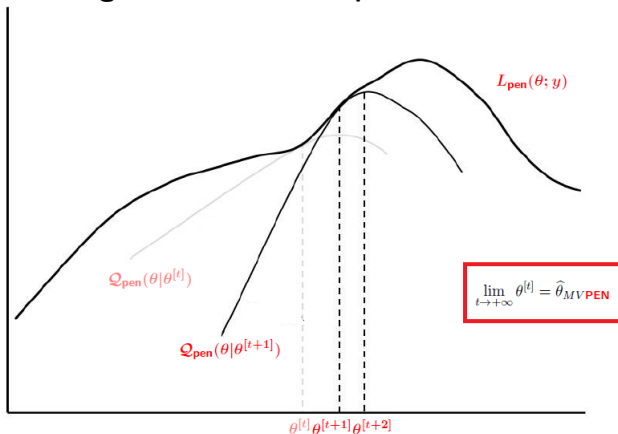
Algorithme EM classique



Algorithme EM avec pénalité RIDGE Green (1990), Eliot et al (2011)



Algorithme EM avec pénalité RIDGE



Estimation par algorithme EM pénalisé

- ▶ Comment calibrer le paramètre de shrinkage λ ?

Estimation par algorithme EM pénalisé

- ▶ Comment calibrer le paramètre de shrinkage λ ?
 - ↪ **Validation Croisée Généralisée (GCV)** à chaque étape de l'EM

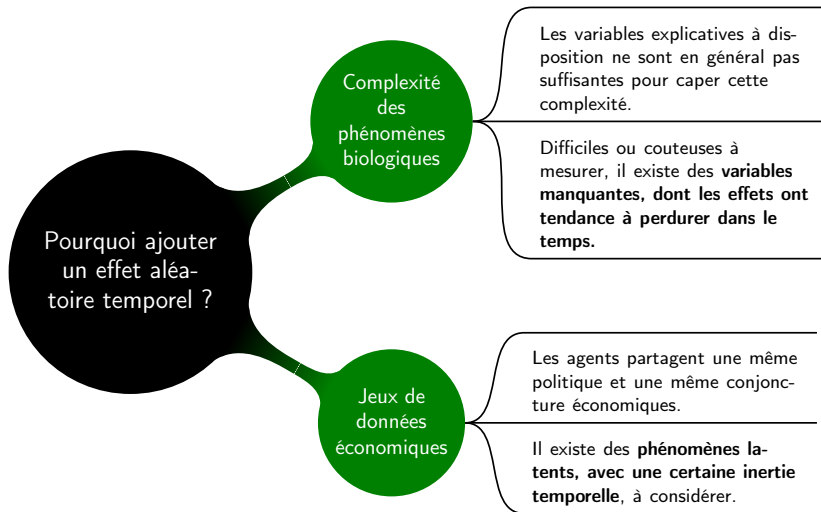
$$\text{GCV} : \lambda^{[t]} \leftarrow \arg \min_{\lambda} \left\{ \text{GCV}(\lambda) = \frac{n^{-1} \left\| y - H_{\lambda}^{[t]} y \right\|^2}{\left[1 - n^{-1} \text{tr} \left(H_{\lambda}^{[t]} \right) \right]^2} \right\}$$

$$\text{E} : Q_{\text{pen}} \left(\theta \mid \theta^{[t]} \right) := \mathbb{E}_{\xi \mid y} \left[L(\theta; y, \xi) - \frac{\lambda^{[t]}}{2} \beta^T \beta \mid \theta^{[t]} \right]$$

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

2 Ridge GL2M pour données longitudinales

- L'existant : L2M-Ridge
- Adaptation aux données longitudinales : ajout d'un effet aléatoire temporel
- Extension aux GL2M
- Simulations et données réelles



L2M pour données longitudinales :

$$y = X\beta + U_1\xi^1 + U_2\xi^2 + \varepsilon$$

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

- ▶ β vecteur des effets fixes
- ▶ $\xi^1 = (\xi_1^1, \xi_2^1, \dots, \xi_N^1)^T$ vecteur de l'effet aléatoire "individu", $U_1 = I_N \otimes \mathbf{1}_R$ matrice de design correspondante
- ▶ $\xi^2 = (\xi_1^2, \xi_2^2, \dots, \xi_R^2)^T$ vecteur de l'effet aléatoire "temporel", $U_2 = \mathbf{1}_N \otimes I_R$ matrice de design correspondante

Modélisation des effets aléatoires

- ▶ $\varepsilon \sim \mathcal{N}(0, \sigma_0^2 I_{NR}) \perp \xi^1 \sim \mathcal{N}(0, \sigma_1^2 I_N)$
- ▶ $\xi^2 \sim \mathcal{N}(0, \Sigma)$, avec $\Sigma_{ij} = \sigma_2^2 \frac{\rho^{|i-j|}}{1-\rho^2}$, $\xi^2 \perp \xi^1$ et $\xi^2 \perp \varepsilon$

Modélisation des effets aléatoires

- ▶ $\varepsilon \sim \mathcal{N}(0, \sigma_0^2 I_{NR}) \perp \xi^1 \sim \mathcal{N}(0, \sigma_1^2 I_N)$
- ▶ $\xi^2 \sim \mathcal{N}(0, \Sigma)$, avec $\Sigma_{ij} = \sigma_2^2 \frac{\rho^{|i-j|}}{1-\rho^2}$, $\xi^2 \perp \xi^1$ et $\xi^2 \perp \varepsilon$

EM Généralisé, avec

$$\theta = (\beta^T, \sigma_0^2, \sigma_1^2, \sigma_2^2, \rho)^T \text{ et } \xi = (\xi^{1T}, \xi^{2T})^T$$

$$\text{GCV} : \lambda^{[t]} \leftarrow \arg \min_{\lambda} \text{GCV}(\lambda)$$

$$\text{E} : Q_{\text{pen}}(\theta | \theta^{[t]}) := \mathbb{E}_{\xi|y} \left[L(\theta; y, \xi) - \frac{\lambda^{[t]}}{2} \beta^T \beta | \theta^{[t]} \right]$$

$$\text{M} : \theta^{[t+1]} \text{ tel que : } Q_{\text{pen}}(\theta^{[t+1]} | \theta^{[t]}) \geq Q_{\text{pen}}(\theta^{[t]} | \theta^{[t]})$$

2 Ridge GL2M pour données longitudinales

- L'existant : L2M-Ridge
- Adaptation aux données longitudinales : ajout d'un effet aléatoire temporel
- **Extension aux GL2M**
- Simulations et données réelles

Extension aux GLMM

$$Y_i | \xi \stackrel{\text{iid}}{\sim} F \text{ de la famille exponentielle}$$
$$g(\underbrace{\mathbb{E}(Y | \xi)}_{\mu}) = \eta = X\beta + U_1\xi^1 + U_2\xi^2$$

Méthodes d'estimation dans les GLMM

- ▶ Approximations d'intégrales : Laplace, Gauss-Hermite
- ▶ Méthodes de Monte Carlo : MCMC, MCEM (McCulloch (1997)), MCML (Knudson (2016))
- ▶ Méthodes de linéarisation : Schall (1991), PQL (Breslow et al. (1993))

La méthode

Étape de LINÉARISATION

- ▶ Linéarisation à l'ordre 1 de y_i au voisinage de μ_i :

$$y_i \simeq z_i = g(\mu_i) + (y_i - \mu_i)g'(\mu_i)$$
$$z_i = \eta_i + e_i$$

- ▶ Définition du modèle linéarisé :

$$\mathcal{M} : z = X\beta + U_1\xi^1 + U_2\xi^2 + e, \quad \text{avec } \mathbb{V}(e) = \Gamma$$

Étape d'ESTIMATION

Algorithme EM pénalisé sur le modèle \mathcal{M}

Itération générique

Linéarisation

$$\mathcal{M} : z^{[t]} = X\beta + U_1\xi^1 + U_2\xi^2 + e, \quad \text{avec } \mathbb{V}(e) = \Gamma^{[t]}$$

Estimation

$$\text{GCV} : \lambda^{[t]} \leftarrow \arg \min_{\lambda} \text{GCV}(\lambda)$$

$$\text{E} : \mathcal{Q}_{\text{pen}}(\theta | \theta^{[t]}) := \mathbb{E}_{\xi|z} \left[L(\theta; z, \xi) - \frac{\lambda^{[t]}}{2} \beta^T \beta | \theta^{[t]} \right]$$

$$\text{M} : \theta^{[t+1]} \text{ tel que : } \mathcal{Q}_{\text{pen}}(\theta^{[t+1]} | \theta^{[t]}) \geq \mathcal{Q}_{\text{pen}}(\theta^{[t]} | \theta^{[t]})$$

Mise à jour

$$\text{Définir } \xi^{[t+1]}, z^{[t+1]}, \Gamma^{[t+1]}$$

2 Ridge GL2M pour données longitudinales

- L'existant : L2M-Ridge
- Adaptation aux données longitudinales : ajout d'un effet aléatoire temporel
- Extension aux GL2M
- Simulations et données réelles

Plan de simulation

Données Poissonniennes

- ▶ $y \sim \mathcal{P}(\exp(X\beta + U_1\xi^1 + U_2\xi^2))$
- ▶ $X = \left[\underbrace{x^1, x^2, x^3, x^4, x^5}_{\text{corr. 2 à 2} = 0.9} \mid \underbrace{x^6, x^7, x^8, x^9, x^{10}}_{\text{indépendantes}} \right]$

Paramètres choisis

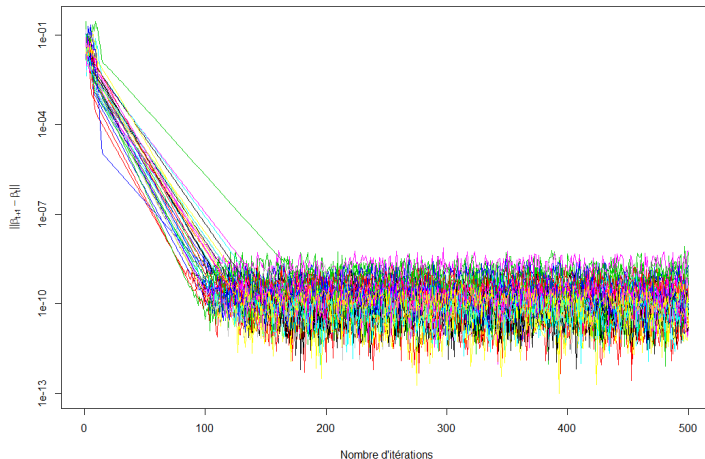
- ▶ $\beta = (0.5, 0.5, 0.5, 0.5, 0.5, 0, 0, 0, 0, 0)$
- ▶ $\sigma_1^2 = \sigma_2^2 = \rho = 0.5$

Individus – Répétitions

- ▶ $N = 10$
- ▶ $R \in \{10, 20, 30, 40, 50, 60, 70, 80, 90, 100\}$

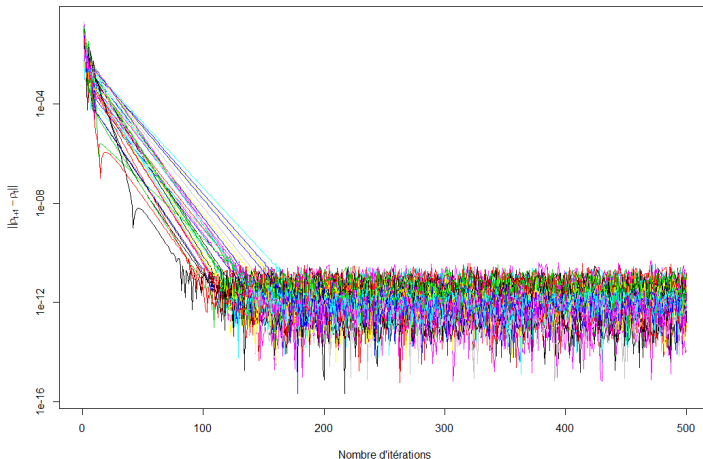
La méthode est-elle convergente ? Effets fixes β

Vitesse de convergence de β



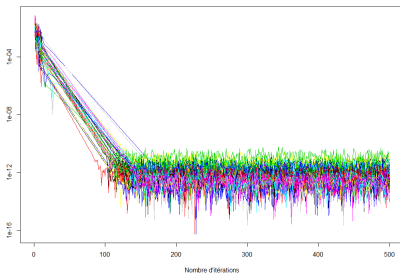
La méthode est-elle convergente ? Paramètre du processus autorégressif ρ

Vitesse de convergence de ρ

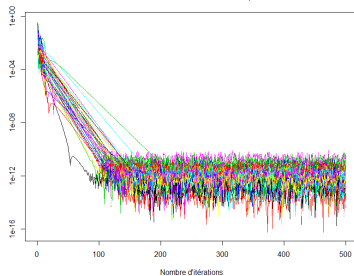


La méthode est-elle convergente ? Paramètres de variances σ_1^2 et σ_2^2

Vitesse de convergence de σ_1^2

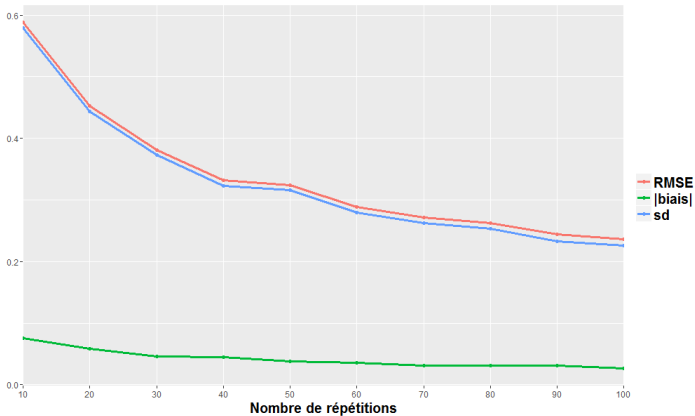


Vitesse de convergence de σ_2^2



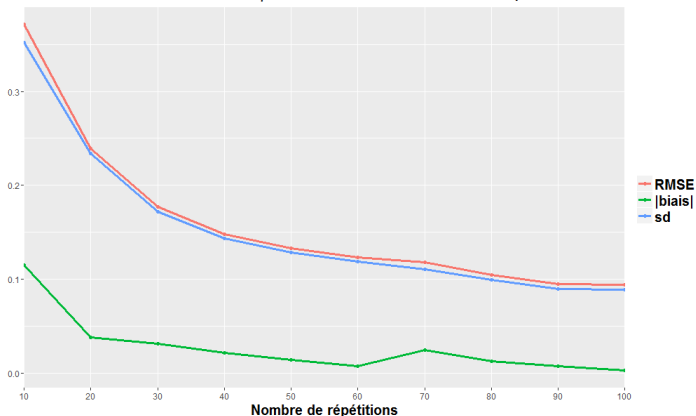
Les paramètres sont-ils bien estimés ? Effets fixes β

Evolution du RMSE de β en fonction du nombre de répétitions



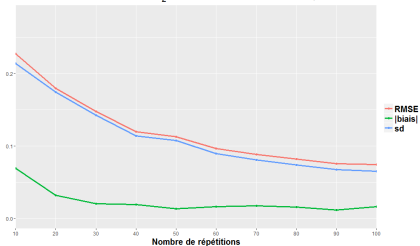
Les paramètres sont-ils bien estimés ? Paramètre du processus autorégressif ρ

Evolution du RMSE de ρ en fonction du nombre de répétitions

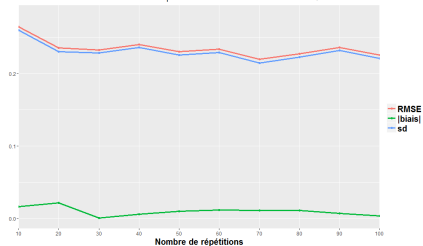


Les paramètres sont-ils bien estimés ? Paramètres de variances σ_1^2 et σ_2^2

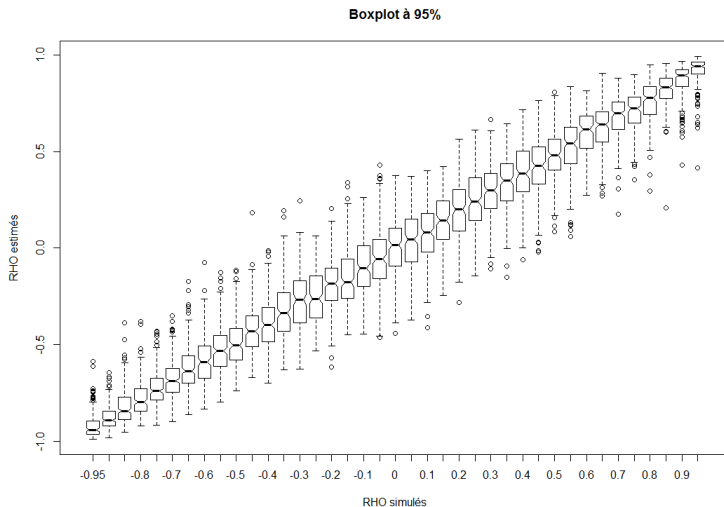
Evolution du RMSE de σ_2^2 en fonction du nombre de répétitions



Evolution du RMSE de σ_1^2 en fonction du nombre de répétitions



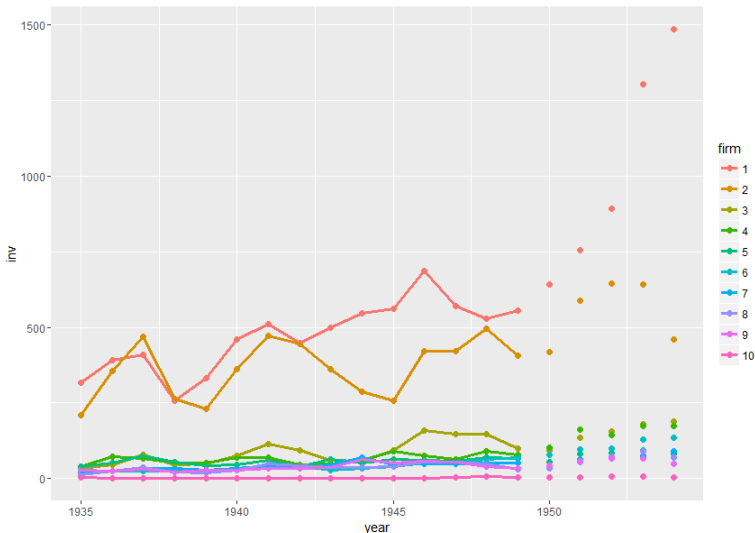
La méthode est-elle sensible à la valeur de ρ ?



Données réelles : "GRUNFELD"

- ▶ Observations annuelles de 10 grandes entreprises des USA sur 20 ans (1935–1954).
- ▶ Variable réponse : "inv" (somme des investissements réalisés par l'entreprise)
- ▶ Variables explicatives : "value" (valeur de l'entreprise sur le marché) et "capital" (capital de l'entreprise).

Modèle calibré sur les 15 premières années



Comparaison des estimations

Package "PLM"

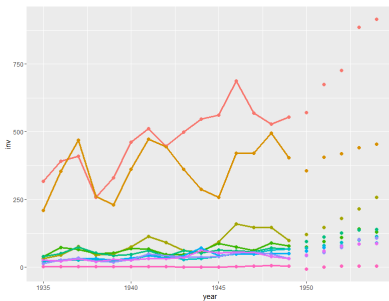
- ▶ $\hat{\beta} = (-4.5, 0.08, 0.20)^T$
- ▶ $\hat{\sigma}_1^2 = 6197$
- ▶ $\hat{\sigma}_2^2 = 95$

Ma méthode

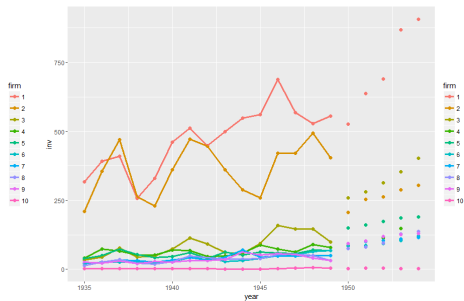
- ▶ $\hat{\beta} = (0.005, 0.07, 0.17)^T$
- ▶ $\hat{\sigma}_1^2 = 5740$
- ▶ $\hat{\sigma}_2^2 = 193$
- ▶ $\hat{\rho} = 0.49$

Comparaison des prédictions

Mes prédictions



vs Leurs prédictions



- 1 Régression Linéaire Généralisée sur Composantes Supervisées (SCGLR)
- 2 Ridge GL2M pour données longitudinales
- 3 Bilan et perspectives

