

Structure Estimation, Graphical Modelling and Causal Inference in High Dimensions: Linear Mixed-effects Models

"**lmmlasso**: Estimation for High-dimensional Linear Mixed-Effects Models Using ℓ_1 -penalization" [5] is a project of C1 in collaboration with C3, and it builds upon [6] and [2].

1. Introduction

n : number of observations, p : number of variables

	$n > p$	$n \ll p$
Linear Regression	Ordinary Least Squares	Lasso
Linear Mixed-Effects Models	Maximum Likelihood (ML) Restricted Maximum Likelihood (REML)	lmmlasso

Key challenges: high-dimensionality, non-convexity

2. Linear mixed models and ℓ_1 -penalized estimation

2.1 High-dimensional Model Set-up

Inhomogeneous data (not independent, but grouped observations)

$i = 1, \dots, N$ grouping index

$j = 1, \dots, n_i$ observation index

$N_T = \sum_{i=1}^N n_i \ll p$ total number of observations

For each group i :

- \mathbf{y}_i : $n_i \times 1$ vector of responses

- \mathbf{X}_i : $n_i \times p$ fixed-effects design matrix

- \mathbf{Z}_i : $n_i \times q$ random-effects design matrix

- \mathbf{b}_i : $q \times 1$ group-specific vector of random regression coefficients

Common for all groups:

- $\boldsymbol{\beta}$: $p \times 1$ vector of fixed regression coefficients

Using the notation from [4], the model can be written as

$$\mathbf{y}_i = \mathbf{X}_i \boldsymbol{\beta} + \mathbf{Z}_i \mathbf{b}_i + \boldsymbol{\varepsilon}_i \quad i = 1, \dots, N, \quad (1)$$

assuming that

i) $\boldsymbol{\varepsilon}_i \sim \mathcal{N}_{n_i}(\mathbf{0}, \sigma^2 \mathbf{I}_{n_i})$ and uncorrelated for $i = 1, \dots, N$,

ii) $\mathbf{b}_i \sim \mathcal{N}_q(\mathbf{0}, \boldsymbol{\Psi})$ and uncorrelated for $i = 1, \dots, N$,

iii) $\boldsymbol{\varepsilon}_1, \dots, \boldsymbol{\varepsilon}_N, \mathbf{b}_1, \dots, \mathbf{b}_N$ are independent.

$\boldsymbol{\Psi} = \boldsymbol{\Psi}_\theta$ is a covariance matrix where θ is an unconstrained set of parameters (with dimension q^*) such that $\boldsymbol{\Psi}_\theta$ is positive definite. From model (1) we deduce that $\mathbf{y}_1, \dots, \mathbf{y}_N$ are independent and $\mathbf{y}_i \sim \mathcal{N}_{n_i}(\mathbf{X}_i \boldsymbol{\beta}, \mathbf{V}_i(\theta, \sigma^2))$ with $\mathbf{V}_i(\theta, \sigma^2) = \mathbf{Z}_i \boldsymbol{\Psi}_\theta \mathbf{Z}_i^T + \sigma^2 \mathbf{I}_{n_i}$.

Denote the stacked vectors $\mathbf{y} = (\mathbf{y}_1^T, \dots, \mathbf{y}_N^T)^T$, $\mathbf{b} = (\mathbf{b}_1^T, \dots, \mathbf{b}_N^T)^T$, $\boldsymbol{\varepsilon} = (\boldsymbol{\varepsilon}_1^T, \dots, \boldsymbol{\varepsilon}_N^T)^T$ and the stacked matrices $\mathbf{X} = (\mathbf{X}_1^T, \dots, \mathbf{X}_N^T)^T$, $\mathbf{Z} = \text{diag}(\mathbf{Z}_1, \dots, \mathbf{Z}_N)$ and $\mathbf{V} = \text{diag}(\mathbf{V}_1, \dots, \mathbf{V}_N)$. Then model (1) can be written as

$$\mathbf{y} = \mathbf{X} \boldsymbol{\beta} + \mathbf{Z} \mathbf{b} + \boldsymbol{\varepsilon} \quad (2)$$

and the negative log-likelihood is given by

$$-\ell(\boldsymbol{\beta}, \boldsymbol{\theta}, \sigma^2) = \frac{1}{2} \left\{ N_T \log(2\pi) + \log |\mathbf{V}| + (\mathbf{y} - \mathbf{X} \boldsymbol{\beta})^T \mathbf{V}^{-1} (\mathbf{y} - \mathbf{X} \boldsymbol{\beta}) \right\} \quad (3)$$

2.2 ℓ_1 -penalized maximum likelihood estimator

Since $N_T \ll p$, consider:

$$Q_\lambda(\boldsymbol{\beta}, \boldsymbol{\theta}, \sigma^2) := \frac{1}{2} \log |\mathbf{V}| + \frac{1}{2} (\mathbf{y} - \mathbf{X} \boldsymbol{\beta})^T \mathbf{V}^{-1} (\mathbf{y} - \mathbf{X} \boldsymbol{\beta}) + \lambda \sum_{k=2}^p |\beta_k|, \quad (4)$$

β_1 : unpenalized intercept, λ : nonnegative regularization parameter

Consequently,

$$\hat{\boldsymbol{\phi}} := (\hat{\boldsymbol{\beta}}, \hat{\boldsymbol{\theta}}, \hat{\sigma}^2) = \underset{\boldsymbol{\beta}, \boldsymbol{\theta}, \sigma^2 > 0, \boldsymbol{\Psi} > 0}{\text{argmin}} Q_\lambda(\boldsymbol{\beta}, \boldsymbol{\theta}, \sigma^2), \quad (5)$$

which is a *nonconvex optimization problem!*

2.3 Prediction of the random effects

Predict the random-effects coefficients \mathbf{b}_i by the *maximum a posteriori (MAP)* principle, which yields

$$\hat{\mathbf{b}}_i = [\mathbf{Z}_i^T \mathbf{Z}_i + \hat{\sigma}^2 \boldsymbol{\Psi}_\theta^{-1}]^{-1} \mathbf{Z}_i^T (\mathbf{y}_i - \mathbf{X}_i \hat{\boldsymbol{\beta}}) \quad i = 1, \dots, N. \quad (6)$$

2.4 Selection of the regularization parameter

Use the Bayesian Information Criterion (BIC) defined by

$$-2\ell(\hat{\boldsymbol{\beta}}, \hat{\boldsymbol{\theta}}, \hat{\sigma}^2) + \log N_T \cdot \hat{d}f_\lambda, \quad (7)$$

where $\hat{d}f_\lambda := |\{1 \leq k \leq p; \hat{\beta}_k \neq 0\}| + \dim(\boldsymbol{\theta})$ is the sum of the number of the nonzero fixed regression coefficients and the number of variance components.

3. Coordinate Gradient Descent Algorithm

We calculate the estimator (5) by the coordinate gradient descent algorithm proposed in [7], and used in [3]. The key elements are:

• **Coordinatewise optimization.** Cycle through the coordinates and minimize the objective function $Q_\lambda(\cdot)$ with respect to only one coordinate while keeping the other parameters fixed.

• **Quadratic approximation.** In each step, approximate $Q_\lambda(\cdot)$ by a strictly convex quadratic function.

• **Inexact line search.** Calculate a descent direction and employ an inexact line search to ensure a decrease in the objective function.

Define

$$P(\boldsymbol{\phi}) := \sum_{k=2}^p |\beta_k|, \quad g(\boldsymbol{\phi}) := \frac{1}{2} \log |\mathbf{V}| + \frac{1}{2} (\mathbf{y} - \mathbf{X} \boldsymbol{\beta})^T \mathbf{V}^{-1} (\mathbf{y} - \mathbf{X} \boldsymbol{\beta}). \quad (8)$$

Then (5) can be written as

$$\hat{\boldsymbol{\phi}} = \underset{\boldsymbol{\phi}}{\text{argmin}} Q_\lambda(\boldsymbol{\phi}) := g(\boldsymbol{\phi}) + \lambda P(\boldsymbol{\phi}) \quad (9)$$

Let \mathbf{e}_j the j th unit vector.

Algorithm. (Coordinate Gradient Descent)

(0) Let $\boldsymbol{\phi}^0 \in \mathbb{R}^{p+q^*+1}$ be an initial value.

For $\ell = 0, 1, 2, \dots$, let \mathcal{S}^ℓ be the index cycling through the coordinates $\{1\}, \{2\}, \dots, \{p+q^*\}, \{p+q^*+1\}$

(1) Approximate the second derivative $\frac{\partial^2}{\partial(\phi_{\mathcal{S}^\ell})^2} Q_\lambda(\boldsymbol{\phi}^\ell)$ by $h^\ell > 0$.

(2) Calculate the descent direction

$$\mathbf{d}^\ell := \underset{\mathbf{d} \in \mathbb{R}}{\text{argmin}} \left\{ g(\boldsymbol{\phi}^\ell) + \frac{\partial}{\partial \phi_{\mathcal{S}^\ell}} g(\boldsymbol{\phi}^\ell) \mathbf{d} + 1/2 d^2 h^\ell + \lambda P(\boldsymbol{\phi}^\ell + \mathbf{d} \mathbf{e}_{\mathcal{S}^\ell}) \right\}.$$

(3) Choose a stepsize $\alpha^\ell > 0$ and set $\boldsymbol{\phi}^{\ell+1} = \boldsymbol{\phi}^\ell + \alpha^\ell \mathbf{d}^\ell \mathbf{e}_{\mathcal{S}^\ell}$ such there is a decrease in the objective function.

until convergence.

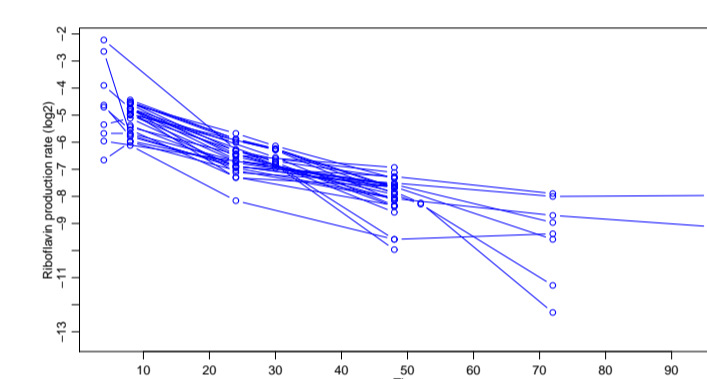
This algorithm is implemented in the R package **lmmlasso**, which is available from the first author's website (<http://stat.ethz.ch/people/schell>) and will be made available on <http://cran.r-project.org>.

4. Application: Riboflavin data

Data description. Data set provided by DSM (Switzerland), see also [1] response: logarithm of the riboflavin production rate of *Bacillus subtilis*

$p = 4088$ covariates measuring the gene expression levels

$N = 28$ groups with $n_i \in \{2, \dots, 6\}$ and $N_T = 111$



Model. We fit a random-intercept model

$$y_{ij} = \mathbf{x}_{ij}^T \boldsymbol{\beta} + b_{i1} + \varepsilon_{ij} \quad i = 1, \dots, N, \quad j = 1, \dots, n_i \quad (10)$$

with $b_{i1} \sim \mathcal{N}(0, \tau^2)$ and $\varepsilon_{ij} \sim \mathcal{N}(0, \sigma^2)$.

Result. Compare **lmmlasso**, *cv-Lasso* (standard Lasso using 10-fold cross-validation) and *Lasso* (standard Lasso using BIC).

	lmmLasso	cv-Lasso	Lasso
$\hat{\sigma}^2$	0.32	0.38	0.30
$\hat{\tau}^2$	0.05	-	-
$ S(\hat{\boldsymbol{\beta}}) $	22	18	21

Conclusions. We see that the total variability can be split into 13.2% *between-subject variability* and 86.8% *within-subject variability*.

5. Current Collaborations

We have three ongoing collaborations:

1) ETH Zurich, Prof. Dr. Sara van de Geer (project C3)

2) ZEW Mannheim, Dr. Stephan Dlugosz, Labour Markets, Human Resources and Social Policy (project A5)

3) TU Berlin, Machine Learning Group, Prof. Dr. Klaus-Robert Müller, Berlin Brain Computer Interface (BBCI)

6. Future Work

In the remaining of the first funding period, we are going to generalize the gaussian linear mixed-effects model to non-gaussian response variables, i.e. the logistic and poisson case. We will focus on theoretical as well as computational aspects and set up an R package called **glmlasso**. This is again joint work with C3.

References

- [1] Kalisch M. Bühlmann, P. and M.H. Maathuis. Variable selection in high-dimensional linear models: partially faithful distributions and the pe-simple algorithm. *Biometrika*, 97:261–278, 2010.
- [2] P. Bühlmann and S. van de Geer. *Statistics for High-Dimensional Data: Methods, Theory and Applications*. to appear, 2011.
- [3] L. Meier, S. van de Geer, and P. Bühlmann. The group lasso for logistic regression. *Journal of the Royal Statistical Society*, 70:53–71, 2008.
- [4] José C. Pinheiro and Douglas M. Bates. *Mixed-Effects Models in S and S-Plus*. Springer, New York, 2000.
- [5] J. Schelldorfer and P. Bühlmann. Estimation for high-dimensional linear mixed-effects models using ℓ_1 -penalization. *arXiv preprint 1002.3784*, 2010.
- [6] N. Städler, P. Bühlmann, and S. van de Geer. ℓ_1 -penalization for mixture regression models (with discussion). *Test*, Online First, 2010.
- [7] P. Tseng and S. Yun. A coordinate gradient descent method for nonsmooth separable minimization. *Mathematical Programming, Series B*, 117:387–423, 2009.



Jörg Schelldorfer
Peter Bühlmann

Research Group FOR916
Statistical Regularization
Project C1

