

# VARIABLE SCREENING AND PARAMETER ESTIMATION FOR HIGH-DIMENSIONAL GENERALIZED LINEAR MIXED MODELS USING $\ell_1$ -PENALIZATION

ETH

Eidgenössische Technische Hochschule Zürich  
Swiss Federal Institute of Technology Zurich

Jürg Schelldorfer and Peter Bühlmann

## 1. Overview

	$n > p$	$n \ll p$
Generalized Linear Models (GLMs)	MLE [glm]	Lasso [glmnet]
Generalized Linear Mixed Models (GLMMs)	MLE [glmer]	<b>GLMMLasso</b> [glmmlasso]

$n$ : number of observations

$p$ : number of variables

## 2. Generalized Linear Mixed Models and $\ell_1$ -Penalized Estimation

### Classical Model Set-up

Notation:

$g = 1, \dots, N$  independent groups/clusters/subjects

$j = 1, \dots, n_g$  observations for group/cluster/subject  $g$

$n = \sum_{g=1}^N n_g$  total number of observations

$\mathbf{y}$ :  $n$ -dim response variable

$\mathbf{b}$ :  $q$ -dim (correlated) random effects

$\beta \in \mathbb{R}^p$  fixed-effects parameters

$\theta \in \mathbb{R}^d$  covariance parameters

$\phi$  dispersion parameter

$\mathbf{X}$ :  $n \times p$  model matrix for  $\beta$

$\mathbf{Z}$ :  $n \times q$  model matrix for  $\mathbf{b}$

$\Sigma_\theta$ :  $q \times q$  covariance matrix, determined by  $\theta$

Model Assumptions (MA 1):

•  $y_i | \mathbf{b}$  are independent for  $i = 1, \dots, n$

•  $y_i | \mathbf{b}$  has a density of the form

$$\exp \left\{ \phi^{-1} \left( y_i \xi_i - b(\xi_i) \right) + c(y_i, \phi) \right\} \text{ with } \mu_i = \mathbb{E}[y_i | \mathbf{b}]$$

•  $g(\mu) = \eta$  with  $\eta = \mathbf{X}\beta + \mathbf{Z}\mathbf{b}$

•  $\mathbf{b} \sim \mathcal{N}_q(\mathbf{0}, \Sigma_\theta)$  with  $\Sigma_\theta \geq 0$  for  $\theta \in \mathbb{R}^d$

Spherical Coordinates (Bates, 2011):

Write  $\Sigma_\theta = \Lambda_\theta \Lambda_\theta^T$  and define  $\mathbf{u}$  by  $\mathbf{b} := \Lambda_\theta \mathbf{u}$  where  $\mathbf{u} \sim \mathcal{N}_q(\mathbf{0}, \mathbf{1}_q)$ .

Parameter Estimation (Bates, 2011):

$$(\hat{\beta}, \hat{\theta}, \hat{\phi})_{MLE} = \arg \min_{\beta, \theta, \phi} -\log L(\beta, \theta, \phi)$$

where  $L(\beta, \theta, \phi)$  is the likelihood function.

### High-dimensional Model Set-up

Model Assumptions (MA 2):

•  $n = \sum_{i=1}^N n_g \ll p$

• the true  $\beta_0$  is sparse

•  $d$  small, say  $d \leq 10$

↔ Goal: Assuming (MA 1) and (MA 2), estimate  $\beta, \theta, \phi$  and predict  $\mathbf{b}$ .

### The GLMMLasso Estimator

To cope with high-dimensionality and to enforce sparsity, we use a **Lasso-type penalty** (Tibshirani, 1996). Hence

$$Q_\lambda(\beta, \theta, \phi) := -2 \log L(\beta, \theta, \phi) + \lambda \|\beta\|_1, \quad (1)$$

for  $\lambda \leq 0$  and

$$(\hat{\beta}, \hat{\theta}, \hat{\phi}) := \arg \min_{\beta, \theta, \phi} Q_\lambda(\beta, \theta, \phi).$$

In general,  $L(\beta, \theta, \phi)$  has no analytical form, so we use the **Laplace approximation** to approximate the integrand of  $L(\beta, \theta, \phi)$  by a quadratic function, i.e.

$$I = \int_{\mathbb{R}^q} e^{-S(\mathbf{u})} d\mathbf{u} \approx \tilde{I}^{LA} = (2\pi)^{q/2} |S''(\tilde{\mathbf{u}})|^{-1/2} e^{-S(\tilde{\mathbf{u}})},$$

where  $\tilde{\mathbf{u}} = \arg \min_{\mathbf{u}} S(\mathbf{u})$  is the mode of  $-S(\mathbf{u})$ .

The Laplace approximation of  $Q_\lambda(\cdot)$  in (1) is

$$\tilde{Q}_\lambda^{LA}(\beta, \theta, \phi) = -2 \sum_{i=1}^n \left\{ \frac{y_i \xi_i(\beta, \theta) - b(\xi_i(\beta, \theta))}{\phi} + c(y_i, \phi) \right\} + \log |(\mathbf{Z}\Lambda_\theta)^T \mathbf{W}_{\beta, \theta, \phi} (\mathbf{Z}\Lambda_\theta) + \mathbf{1}_q| + \|\tilde{\mathbf{u}}(\beta, \theta, \phi)\|_2^2 + \lambda \|\beta\|_1,$$

where  $\mathbf{W}_{\beta, \theta, \phi} = \text{diag}^{-1} \left( \phi v(\mu_i(\beta, \theta)) g'(\mu_i(\beta, \theta))^2 \right)_{i=1}^n$  and  $v(\cdot)$  is the variance function.

The **GLMMLasso** estimator is defined by

$$(\hat{\beta}, \hat{\theta}, \hat{\phi}) := \arg \min_{\beta, \theta, \phi} \tilde{Q}_\lambda^{LA}(\beta, \theta, \phi) \quad (2)$$

This is a high-dimensional, non-convex optimization problem!

## 3. Computational Algorithm

To calculate the GLMMLasso estimator (2), we use **coordinatewise optimization with inexact line search**, i.e. optimizing  $\tilde{Q}_\lambda^{LA}(\cdot)$  with respect to one coordinate keeping all other coordinates fixed (Tseng and Yun, 2009). An overview of the algorithm may be described as follows:

### GLMMLasso algorithm

(0) Choose a starting value  $(\beta^{(0)}, \theta^{(0)}, \phi^{(0)})$ .

**Repeat** for  $s = 1, 2, \dots$

(1) (FIXED-EFFECTS PARAMETER OPTIMIZATION)

For  $k = 1, \dots, p$

a) (Laplace approximation)

Calculate the Laplace approximation  $\tilde{Q}_\lambda^{LA}(\cdot)$  based on the current parameter estimates.

b) (Quadratic approximation and inexact line search)

i) Approximate the second derivative by  $h_k^{(s)} > 0$ .

ii) Calculate the descent direction  $d_k^{(s)} \in \mathbb{R}$

iii) Choose a step size  $\alpha_k^{(s)} > 0$  such that there is a decrease in the objective function.

(2) (COVARIANCE PARAMETER OPTIMIZATION)

For  $l = 1, \dots, d$

$$\theta_l^{(s)} = \arg \min_{\theta_l} \tilde{Q}_\lambda^{LA}(\cdot)$$

(3) (DISPERSION PARAMETER OPTIMIZATION)

$$\phi^{(s)} = \arg \min_{\phi} \tilde{Q}_\lambda^{LA}(\cdot)$$

**until** convergence.

The above algorithm solves (2) exactly. In order to speed up the algorithm, we suggest an **approximate algorithm** comprising the following two parts:

• Regard  $\tilde{\mathbf{u}}$  as fixed for the quadratic approximation with respect to the derivatives of  $\beta_k$  in (1) b).

• Active-set algorithm: cycle through the non-zero coefficients  $\beta_k^{(s)}$ , and through all  $p$  coefficients only every  $D$ th iteration.

## 4. Two-stage GLMMLasso estimator(s)

Apart from good variable selection properties accomplished by the Lasso, we advocate a two-stage procedure to get accurate parameter estimates. Hence the first stage aims at estimating a candidate set of variables (*variable screening*). The goal of the second step is unbiased parameter estimation (*parameter estimation*). Therefore we propose a **refitting by ML methods**. This two-stage procedure can be summarized as follows:

### Two-stage GLMMLasso

Stage 1: Compute the GLMMLasso estimator (2).

Stage 2: Perform a ML method as in (3) or (4).

Let  $(\hat{\beta}_{init}, \hat{\theta}_{init}, \hat{\phi}_{init})$  denote the estimate from (2).

### The GLMMLasso-MLE hybrid estimator

Define  $\hat{S}_{hybrid} := \{k : |\hat{\beta}_{k,init}| \neq 0\}$ . Then

$$(\hat{\beta}, \hat{\theta}, \hat{\phi})_{hybrid} := \arg \min_{\beta_{\hat{S}_{hybrid}}, \theta, \phi} -2 \log L(\beta_{\hat{S}_{hybrid}}, \theta, \phi) \quad (3)$$

where for  $S \subseteq \{1, \dots, p\}$ ,  $(\beta_S)_k = \beta_k$  if  $k \in S$  and  $(\beta_S)_k = 0$  if  $k \notin S$ .

### The thresholded GLMMLasso estimator

The thresholded Lasso with refitting was examined in Geer et al (2010). Define  $\hat{S}_{thres} := \{k : |\hat{\beta}_{k,init}| > \lambda_{thres}\}$ . The thresholded GLMMLasso estimator is defined by

$$(\hat{\beta}, \hat{\theta}, \hat{\phi})_{thres} := \arg \min_{\beta_{\hat{S}_{thres}}, \theta, \phi} -2 \log L(\beta_{\hat{S}_{thres}}, \theta, \phi) \quad (4)$$

### Selection of the regularization parameters

For the choice of the regularization parameters  $\lambda$  and  $\lambda_{thres}$ , we propose the BIC and the AIC

$$c_{n,\lambda} = -2 \log L(\hat{\beta}, \hat{\theta}, \hat{\phi}) + a(n) \cdot \hat{d}f_\lambda$$

where  $a(n) = \log(n)$  for the BIC,  $a(n) = 2$  for the AIC and  $\hat{d}f_\lambda = |\{1 \leq k \leq p : \hat{\beta}_k \neq 0\}| + \dim(\hat{\theta})$ .

## 6. Illustration

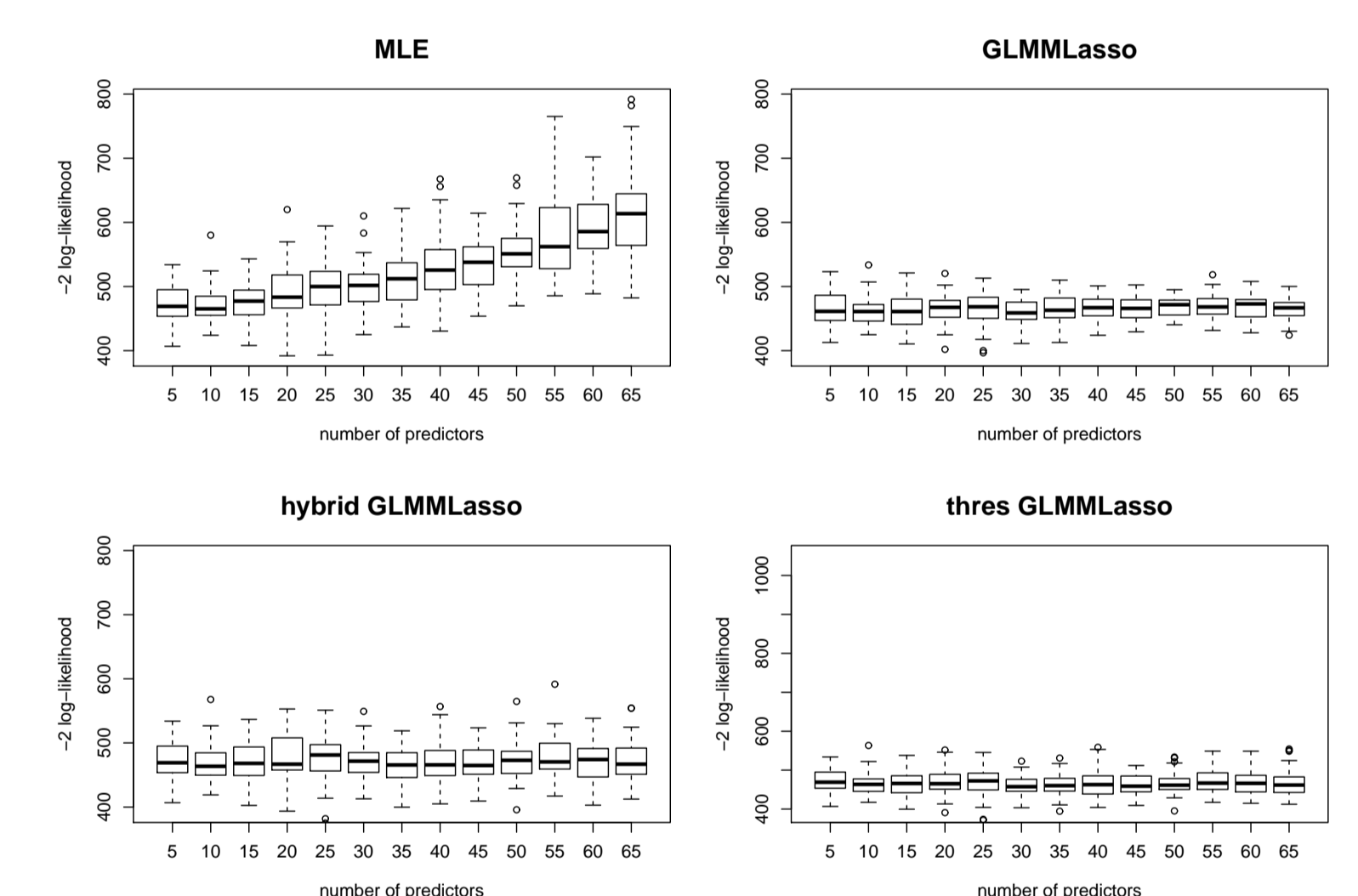


FIGURE 1: Minus twice out-of-sample log-likelihood for a growing number of covariates. The MLE performs badly whereas the GLMMLasso estimators remain stable. We use a random-intercept logistic mixed model with  $n = 400$ ,  $N = 40$ ,  $n_g = 10$ ,  $\theta^2 = 1$ ,  $\beta_0 = (0, 1, -1, 1, -1)$ .

## 6. R package

An implementation of the algorithm will be available online in the R package **glmmlasso**. The Gaussian case is implemented in the standalone R package **lmmlasso** (Schelldorfer et al, 2011), which is available from CRAN.

## References

- Bates D. M. (2011) Computational Methods for Mixed Models. available at <http://cran.r-project.org/web/packages/lme4/vignettes/Theory.pdf>
- Tibshirani R. (1996) Regression Shrinkage and Selection via the Lasso. Journal of the Royal Statistical Society, Series B 58:267–288
- Tseng P. and Yun S. (2009) A Coordinate Gradient Descent Method for Nonsmooth Separable Minimization. Mathematical Programming 117:387–423
- Schelldorfer J., Bühlmann P., van de Geer S. (2011) Estimation for High-Dimensional Linear Mixed-Effects Models Using  $\ell_1$ -penalization. The Scandinavian Journal of Statistics 38:197–214
- van de Geer S., Bühlmann P., Zhou S. (2010) The Adaptive and the Thresholded Lasso for Potentially Misspecified Models. arXiv:1001.5176v3