

GLMMLasso: An Algorithm for High-Dimensional Generalized Linear Mixed Models Using ℓ_1 -Penalization

Jürg SCHELLDORFER, Lukas MEIER, and Peter BÜHLMANN

We propose an ℓ_1 -penalized algorithm for fitting high-dimensional generalized linear mixed models (GLMMs). GLMMs can be viewed as an extension of generalized linear models for clustered observations. Our Lasso-type approach for GLMMs should be mainly used as variable screening method to reduce the number of variables below the sample size. We then suggest a refitting by maximum likelihood based on the selected variables only. This is an effective correction to overcome problems stemming from the variable screening procedure that are more severe with GLMMs than for generalized linear models. We illustrate the performance of our algorithm on simulated as well as on real data examples. Supplementary materials are available online and the algorithm is implemented in the R package `glmixedlasso`.

Key Words: Coordinate gradient descent; Laplace approximation; Random-effects model; Variable selection.

1. INTRODUCTION

In recent years, high-dimensional linear regression models have been extensively studied. The most popular method to achieve sparse estimates is the Lasso (Tibshirani 1996), which uses an ℓ_1 -penalty. The Lasso is attractive not only in terms of its statistical properties but also due to its fast computation solving a convex optimization problem. However, relatively few articles examining high-dimensional regression problems involving a nonconvex loss function can be cited—for example, Khalili and Chen (2007) and Städler, Bühlmann, and van de Geer (2010) for Gaussian mixture models; Pan and Shen (2007) and Witten and Tibshirani (2010) for clustering; and Witten and Tibshirani (2011) for linear discriminant analysis.

Generalized linear mixed models or GLMMs (McCullagh and Nelder 1989; Breslow and Clayton 1993; McCulloch and Searle 2001; Molenberghs and Verbeke 2005) are an

Jürg Schelldorfer is Statistician, AXA Winterthur, CH-8400 Winterthur, Switzerland (E-mail: juerg.schelldorfer@axa-winterthur.ch). Lukas Meier is Senior Scientist (E-mail: meier@stat.math.ethz.ch) and Peter Bühlmann is Professor (E-mail: buhlmann@stat.math.ethz.ch), Seminar für Statistik, ETH Zürich, CH-8092 Zürich, Switzerland.

extension of generalized linear models by adding random effects to the linear predictor to accommodate for clustered or overdispersed data. These models have received much attention in many applications such as biology, ecology, medicine, pharmaceutical science, and econometrics. Available software packages (`lme4` in R, NLMIXED in SAS, among others) allow to fit a wide range of GLMMs.

In this article, we develop a method for high-dimensional GLMMs. It is based on a Lasso-type regularization with a cyclic coordinate descent optimization. Due to shrinkage introduced by ℓ_1 -penalization, our approach performs in a first step variable screening, thereby selecting a set of candidate active variables. In other words, the proposed method primarily aims at reducing the dimensionality of the high-dimensional GLMM. In a second step, we perform refitting by maximum likelihood (ML) estimation to get accurate parameter estimates. The idea of such a two-stage approach has been used in linear models (Efron et al. 2004) and it is related to the adaptive Lasso (Zou 2006) and the thresholded Lasso (Zhou 2010; van de Geer, Bühlmann, and Zhou 2011). In fact, a two-stage approach is much more important than for linear models since shrinkage in GLMMs can have a severe effect on the estimation of variance components (see Sections 4 and 5).

To the best of our knowledge, there does not exist any literature devoted to truly high-dimensional GLMMs. Some papers focus on penalized variable selection procedures in generalized mixed models with low-dimensional data: we refer to Yang (2007), Ibrahim et al. (2010), and Ni, Zhang, and Zhang (2010). Groll and Tutz (in press) have independently studied the same statistical problem and have also used a Lasso-type approach but with a focus on rather low-dimensional problems. Few papers focus on variable selection in generalized additive mixed models (e.g., Xue, Qu, and Zhou 2010; Lai, Huang, and Lee 2012). Schelldorfer, Bühlmann, and van de Geer (2011) presented statistical theory and an algorithm for high-dimensional Gaussian linear mixed models, where computation is much easier than in the generalized case.

The main contribution of the present article is the construction and implementation of an efficient algorithm for ℓ_1 -penalization in truly high-dimensional GLMMs, called the GLMMLasso. We use the Laplace approximation (Bates 2011b) and combine it with efficient coordinate gradient descent (CGD) methods (Tseng and Yun 2009). Our algorithm is feasible for problems where the number of variables is in the thousands and taking advantage of sparsity with respect to dimensionality (i.e., only few active variables) is exploited by an active set strategy.

The rest of the article is organized as follows. In Section 2, we review the GLMM and introduce the GLMMLasso estimator. In Section 3, we describe the details of the computational algorithm before advocating the two-stage GLMMLasso estimators in Section 4. In Sections 5 and 6, we consider the performance of our methods on simulated and real datasets. The article concludes with a discussion in Section 7. Supplementary materials including additional simulation examples are available online.

2. GENERALIZED LINEAR MIXED MODELS AND ℓ_1 -PENALIZED ESTIMATION

In this section, we first look at the classical GLMM setting where the number of observations is larger than the number of covariates, that is, $p < n$. We closely follow

Bates (2011a). Second, we consider the high-dimensional framework, that is, $n \ll p$, and present the ℓ_1 -penalized ML estimator.

2.1 MODEL FORMULATION

Suppose that the observations are not independent but grouped instead. Let $r = 1, \dots, N$ be the grouping index and $j = 1, \dots, n_r$ the j th outcome within group r . Denote by n the total number of observations, that is, $n = \sum_{r=1}^N n_r$. Let \mathbf{X} be the $n \times p$ fixed-effects design matrix, \mathbf{Z} the $n \times q$ random-effects design matrix, \mathbf{Y} the n -dimensional random response vector, and \mathbf{B} be the q -dimensional vector of random effects. We observe \mathbf{y} of \mathbf{Y} whereas \mathbf{B} is unobserved. The GLMM is specified by the unconditional distribution of \mathbf{B} and the conditional distribution of $\mathbf{Y}|\mathbf{B} = \mathbf{b}$:

- (i) $\mathcal{Y}_i|\mathbf{B} = \mathbf{b}$ are independent for $i = 1, \dots, n$.
- (ii) The distribution of $\mathcal{Y}_i|\mathbf{B} = \mathbf{b}$ belongs to the exponential family with density

$$\exp\{\phi^{-1}(y_i \xi_i - b(\xi_i)) + c(y_i, \phi)\},$$

where $b(\cdot)$ and $c(\cdot, \cdot)$ are known functions. ϕ is the dispersion parameter (known or unknown) and ξ_i is associated with the conditional mean $\mu_i := E[\mathcal{Y}_i|\mathbf{B} = \mathbf{b}]$, that is, $\xi_i = \xi_i(\mu_i)$.

- (iii) The conditional mean vector $\boldsymbol{\mu}$ depends on \mathbf{b} through the known link function g and the linear predictor $\boldsymbol{\eta} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{b}$, with $\boldsymbol{\eta} = \mathbf{g}(\boldsymbol{\mu})$ componentwise. Here, $\boldsymbol{\beta}$ is the unknown p -dimensional parameter vector, called fixed effects, and \mathbf{b} the unknown q -dimensional vector of random effects.
- (iv) $\mathbf{B} \sim \mathcal{N}_q(\mathbf{0}, \boldsymbol{\Sigma}_\theta)$ where the covariance matrix $\boldsymbol{\Sigma}_\theta$ is parameterized by the unknown parameter vector $\boldsymbol{\theta} \in \mathbb{R}^d$. We assume that $\boldsymbol{\Sigma}_\theta$ is positive semidefinite, that is, $\boldsymbol{\Sigma}_\theta \geq 0$. The dimensionality d is typically small, say $d \leq 10$.

By using \mathbf{B} and $\boldsymbol{\Sigma}_\theta$ in the definition above, we have already defined the random-effects structure of the GLMM. To be more precise, we have specified which variables have an additional random effect and how the structure of $\boldsymbol{\Sigma}_\theta$ looks like (e.g., multiple of the identity or diagonal). A discussion of how to find these structures is beyond the scope of this article.

Let us write $\boldsymbol{\Sigma}_\theta$ in terms of its Cholesky decomposition $\boldsymbol{\Sigma}_\theta = \boldsymbol{\Lambda}_\theta \boldsymbol{\Lambda}_\theta^T$ and introduce the (unobserved) random variable \mathbf{U} defined by $\mathbf{B} := \boldsymbol{\Lambda}_\theta \mathbf{U}$ where $\mathbf{U} \sim \mathcal{N}_q(\mathbf{0}, \mathbf{1}_q)$. Then the linear predictor $\boldsymbol{\eta}$ can be written as $\boldsymbol{\eta} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\boldsymbol{\Lambda}_\theta \mathbf{u}$. We estimate the parameters $\boldsymbol{\beta}$, $\boldsymbol{\theta}$, and ϕ (if unknown) by the ML method and predict the random effects \mathbf{u} .

2.2 LIKELIHOOD FUNCTION

Employing the notation $\xi_i(\mu_i) = \xi_i(\boldsymbol{\beta}, \boldsymbol{\theta})$, the likelihood function of a GLMM is given by the following expression:

$$L(\boldsymbol{\beta}, \boldsymbol{\theta}, \phi) = \int_{\mathbb{R}^q} \prod_{i=1}^n [\exp\{\phi^{-1}(y_i \xi_i(\boldsymbol{\beta}, \boldsymbol{\theta}) - b(\xi_i(\boldsymbol{\beta}, \boldsymbol{\theta}))) + c(y_i, \phi)\}] \\ \times \frac{1}{(2\pi)^{q/2}} \exp\left\{-\frac{1}{2}\|\mathbf{u}\|_2^2\right\} d\mathbf{u}$$

$$= \frac{1}{(2\pi)^{q/2}} \int_{\mathbb{R}^q} \exp \left\{ \sum_{i=1}^n \left(\frac{y_i \xi_i(\boldsymbol{\beta}, \boldsymbol{\theta}) - b(\xi_i(\boldsymbol{\beta}, \boldsymbol{\theta}))}{\phi} + c(y_i, \phi) \right) - \frac{1}{2} \|\mathbf{u}\|_2^2 \right\} d\mathbf{u}. \quad (1)$$

In general, the integral (1) cannot be worked out analytically and numerical approximations are required (see Skrondal and Rabe-Hesketh 2004; Molenberghs and Verbeke 2005; Jiang 2007).

2.3 THE GLMMLASSO ESTIMATOR

We now turn to the high-dimensional setting where the number of fixed-effect variables p is much larger than the number of observations n , that is, we study the so-called $n \ll p$ setup.

Let us assume that the true underlying fixed-effects vector $\boldsymbol{\beta}_0$ is sparse in the sense that many coefficients of $\boldsymbol{\beta}_0$ are zero. To enforce sparsity of our estimator, we advocate a Lasso-type approach. This means that we add an ℓ_1 -penalty for the fixed-effects vector $\boldsymbol{\beta}$ to the likelihood function. Thus, we are going to consider the following objective function:

$$Q_\lambda(\boldsymbol{\beta}, \boldsymbol{\theta}, \phi) = -2 \log L(\boldsymbol{\beta}, \boldsymbol{\theta}, \phi) + \lambda \|\boldsymbol{\beta}\|_1, \quad (2)$$

where $\lambda \geq 0$ is a regularization parameter. Appropriate choices for λ are discussed in Section 4.

We aim at estimating the fixed-effect parameter $\boldsymbol{\beta}$, the covariance parameter $\boldsymbol{\theta}$, and, if unknown, the dispersion parameter ϕ , by

$$(\hat{\boldsymbol{\beta}}, \hat{\boldsymbol{\theta}}, \hat{\phi}) := \arg \min_{\boldsymbol{\beta}, \boldsymbol{\theta}, \phi} Q_\lambda(\boldsymbol{\beta}, \boldsymbol{\theta}, \phi). \quad (3)$$

We call (3) the GLMMLasso estimator. Since the likelihood function (1) comprises analytically intractable integrals (except for the Gaussian case), some approximations have to be used. We are going to illustrate the algorithm using the Laplace approximation. For GLMMs, it is accurate with low computational burden, as advocated by Bates (2011b). A thorough discussion of the accuracy and limitations of the Laplace approximation can be found in the article by Joe (2008). Generally, the Laplace approximation is used to calculate integrals of the form

$$I = \int_{\mathbb{R}^q} e^{-S(\mathbf{u})} d\mathbf{u}, \quad (4)$$

where $S(\mathbf{u})$ is a known function of a q -dimensional variable \mathbf{u} . Let

$$\tilde{\mathbf{u}} = \arg \max_{\mathbf{u}} -S(\mathbf{u}) \quad (5)$$

(i.e., $S'(\tilde{\mathbf{u}}) = 0$), then the Laplace approximation of I is given by

$$I \approx I^{LA} = (2\pi)^{q/2} |S''(\tilde{\mathbf{u}})|^{-1/2} e^{-S(\tilde{\mathbf{u}})}. \quad (6)$$

The mode $\tilde{\mathbf{u}}$ in (5) is calculated by the penalized iterative least squares (PIRLS) algorithm. It is presented in the literature by Bates (2011b) and described in the online supplementary materials. The PIRLS algorithm is related to the iterative reweighted least squares (IRLS) algorithm for obtaining the ML estimator in generalized linear models.

It should be noted that $\tilde{\mathbf{u}}$ depends on $\boldsymbol{\beta}$, $\boldsymbol{\theta}$, and ϕ . From (1) and (6), we deduce that the Laplace approximation of the objective function $Q_\lambda(\cdot)$ in (2) is

$$Q_\lambda^{\text{LA}}(\boldsymbol{\beta}, \boldsymbol{\theta}, \phi) = -2 \sum_{i=1}^n \left\{ \frac{y_i \xi_i(\boldsymbol{\beta}, \boldsymbol{\theta}) - b(\xi_i(\boldsymbol{\beta}, \boldsymbol{\theta}))}{\phi} + c(y_i, \phi) \right\} \\ + \log |(\mathbf{Z}\boldsymbol{\Lambda}_\theta)^T \mathbf{W}_{\boldsymbol{\beta}, \boldsymbol{\theta}, \phi}(\mathbf{Z}\boldsymbol{\Lambda}_\theta) + \mathbf{1}_q| + \|\tilde{\mathbf{u}}(\boldsymbol{\beta}, \boldsymbol{\theta}, \phi)\|_2^2 + \lambda \|\boldsymbol{\beta}\|_1, \quad (7)$$

where $\mathbf{W}_{\boldsymbol{\beta}, \boldsymbol{\theta}, \phi} = \text{diag}^{-1}(\phi v(\mu_i(\boldsymbol{\beta}, \boldsymbol{\theta})) g'(\mu_i(\boldsymbol{\beta}, \boldsymbol{\theta}))^2)_{i=1}^n$ and $v(\cdot)$ is the known conditional variance function (McCullagh and Nelder 1989). The estimator (3) is then approximated by

$$(\hat{\boldsymbol{\beta}}^{\text{LA}}, \hat{\boldsymbol{\theta}}^{\text{LA}}, \hat{\phi}^{\text{LA}}) := \arg \min_{\boldsymbol{\beta}, \boldsymbol{\theta}, \phi} Q_\lambda^{\text{LA}}(\boldsymbol{\beta}, \boldsymbol{\theta}, \phi). \quad (8)$$

We call (8) the GLMMLasso^{LA} estimator. It is the approximation (8) to the objective function (3) that is optimized to obtain the parameter estimates. Moreover, we would like to emphasize that (8) is a nonconvex function with respect to $(\boldsymbol{\beta}, \boldsymbol{\theta}, \phi)$ consisting of a nonconvex loss function and a convex penalty.

3. COMPUTATIONAL ALGORITHM

In this section, we present the computational algorithm to obtain the GLMMLasso^{LA} estimator (8). The algorithm is based on ideas by Tseng and Yun (2009) of the (block) CGD method. The notion of the CGD algorithm is that we cycle through components of the full parameter vector $\boldsymbol{\psi} := (\boldsymbol{\beta}, \boldsymbol{\theta}, \phi) \in \mathbb{R}^{p+d+1}$ and minimize the objective function $Q_\lambda^{\text{LA}}(\cdot)$ only with respect to one parameter while keeping the other parameters fixed. In doing so, we calculate a quadratic approximation and perform an indirect line search to ensure that the objective function decreases. (Block) CGD algorithms are used by Meier, van de Geer, and Bühlmann (2008), Wu and Lange (2008), Friedman, Hastie, and Tibshirani (2010), and Breheny and Huang (2011), and are now extremely popular in high-dimensional penalized regression problems.

We first give an overview of the algorithm that solves minimization problem (8) exactly before considering an approximate algorithm that finds a solution close to the exact minimizer of (8). Finally, we present some details of the algorithm.

3.1 THE EXACT GLMMLASSO ALGORITHM

We describe here an exact algorithm, called exact GLMMLasso (we notationally omit the involved Laplace approximation), for the Laplace approximated objective function in (8). Let us write (7) with a different notation to ease the presentation. For $\boldsymbol{\psi} = (\boldsymbol{\beta}, \boldsymbol{\theta}, \phi) \in \mathbb{R}^{p+d+1}$, define the function

$$f(\boldsymbol{\psi}) := -2 \sum_{i=1}^n \left\{ \frac{y_i \xi_i(\boldsymbol{\beta}, \boldsymbol{\theta}) - b(\xi_i(\boldsymbol{\beta}, \boldsymbol{\theta}))}{\phi} + c(y_i, \phi) \right\} \\ + \log |(\mathbf{Z}\boldsymbol{\Lambda}_\theta)^T \mathbf{W}_\psi(\mathbf{Z}\boldsymbol{\Lambda}_\theta) + \mathbf{1}_q| + \|\tilde{\mathbf{u}}(\boldsymbol{\psi})\|_2^2.$$

Now (8) can be written as $\hat{\boldsymbol{\psi}}_{\lambda}^{\text{LA}} = \arg \min_{\boldsymbol{\psi}} Q_{\lambda}^{\text{LA}}(\boldsymbol{\psi}) := f(\boldsymbol{\psi}) + \lambda \|\boldsymbol{\beta}\|_1$. Let \mathbf{e}_j be the j th unit vector and denote by (s) the s th iteration step. Moreover, we let

$$\boldsymbol{\beta}^{(s)} := (\beta_1^{(s)}, \dots, \beta_p^{(s)})^T, \quad \boldsymbol{\theta}^{(s)} := (\theta_1^{(s)}, \dots, \theta_d^{(s)})^T, \quad \phi^{(s)}$$

be the estimates of $\boldsymbol{\beta}$, $\boldsymbol{\theta}$, and ϕ in the s th iteration. Using the notation

$$\begin{aligned} \boldsymbol{\beta}^{(s,s-1;\beta_k)} &:= (\beta_1^{(s)}, \dots, \beta_{k-1}^{(s)}, \beta_k, \beta_{k+1}^{(s-1)}, \dots, \beta_p^{(s-1)})^T, \\ \boldsymbol{\theta}^{(s,s-1;\theta_l)} &:= (\theta_1^{(s)}, \dots, \theta_{l-1}^{(s)}, \theta_l, \theta_{l+1}^{(s-1)}, \dots, \theta_d^{(s-1)})^T, \\ \boldsymbol{\beta}^{(s,s-1;k)} &:= (\beta_1^{(s)}, \dots, \beta_{k-1}^{(s)}, \beta_k^{(s-1)}, \beta_{k+1}^{(s-1)}, \dots, \beta_p^{(s-1)})^T, \end{aligned}$$

the exact GLMMLasso algorithm is summarized in Algorithm 1.

Particularly in the high-dimensional setting, the calculation of the quadratic approximation requires a large amount of computing time. Therefore, it is interesting to examine a much faster approximate algorithm.

3.2 THE (APPROXIMATE) GLMMLASSO ALGORITHM

In the exact Algorithm 1, we consider in Step 1b the mode $\tilde{\boldsymbol{u}}$ as a function of the parameters, that is, $\tilde{\boldsymbol{u}} = \tilde{\boldsymbol{u}}(\boldsymbol{\beta}, \boldsymbol{\theta}, \phi)$. However, the calculation of the derivatives of $f(\cdot)$ with respect to β_k is computationally intensive. This becomes a major issue in the high-dimensional setting where a substantial amount of computing time is allocated to this particular part of the algorithm. In addition, the exact GLMMLasso algorithm requires a large number of outer iterations s . To attenuate these difficulties, we propose a slightly modified version of Algorithm 1. We suggest performing the quadratic approximation and the inexact line search while considering $\tilde{\boldsymbol{u}}$ as fixed and not depending on β_k . Denoting by $f(\cdot|\tilde{\boldsymbol{u}})$ the function $f(\cdot)$ for which $\tilde{\boldsymbol{u}}$ is considered as fixed, the (approximate) GLMMLasso algorithm is given in Algorithm 2:

We illustrate in the online supplementary materials that the approximate GLMMLasso algorithm speeds up remarkably without losing that much accuracy. Additionally, the approximation emphasizes the importance of a refitting as advocated in the next section.

3.3 CONVERGENCE BEHAVIOR AND DETAILS OF THE GLMMLASSO ALGORITHM

3.3.1 Numerical Convergence. The convergence of the exact GLMMLasso algorithm to a stationary point can be proofed using the results presented by Tseng and Yun (2009). It is worth pointing out that in the low-dimensional framework, the exact GLMMLasso algorithm with $\lambda = 0$ (no penalization) gives the same results as the function `glmmer` in the R package `lme4`.

- (0) *Starting value $\boldsymbol{\psi}^{(0)}$.* As starting value for $\boldsymbol{\beta}$, we fit a generalized linear model with the Lasso where the regularization parameter is chosen by cross-validation. The initial values for $\boldsymbol{\theta}$ and ϕ are then calculated using Steps (2) and (3) in Algorithms 1 and 2.

Algorithm 1 *Exact GLMMLasso algorithm*

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

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

(1) *(fixed-effect parameter optimization)*

For $k = 1, \dots, p$

a) *(Laplace approximation)*

Calculate the Laplace approximation

$$Q_{\lambda}^{\text{LA}}(\boldsymbol{\beta}^{(s,s-1;k)}, \boldsymbol{\theta}^{(s-1)}, \phi^{(s-1)}).$$

b) *(Quadratic approximation and inexact line search)*

i) Approximate the second derivative

$$\frac{\partial^2}{\partial \beta_k^2} f(\boldsymbol{\beta}^{(s,s-1;\beta_k)}, \boldsymbol{\theta}^{(s-1)}, \phi^{(s-1)}) \Big|_{\beta_k = \beta_k^{(s-1)}}$$

by $h_k^{(s)} > 0$ as described in the section below.

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

$$d_k^{(s)} := \arg \min_d \left\{ f(\boldsymbol{\beta}^{(s,s-1;k)}, \boldsymbol{\theta}^{(s-1)}, \phi^{(s-1)}) + \frac{\partial}{\partial \beta_k} f(\boldsymbol{\beta}^{(s,s-1;\beta_k)}, \boldsymbol{\theta}^{(s-1)}, \phi^{(s-1)}) \Big|_{\beta_k = \beta_k^{(s-1)}} d + \frac{1}{2} d^2 h_k^{(s)} + \lambda \|\boldsymbol{\beta}^{(s,s-1;k)} + d \mathbf{e}_k\|_1 \right\}.$$

iii) Choose a step size $\alpha_k^{(s)} > 0$ and set $\boldsymbol{\beta}^{(s,s-1;k+1)} = \boldsymbol{\beta}^{(s,s-1;k)} + \alpha_k^{(s)} d_k^{(s)} \mathbf{e}_k$ such that

$$Q_{\lambda}^{\text{LA}}(\boldsymbol{\beta}^{(s,s-1;k+1)}, \boldsymbol{\theta}^{(s-1)}, \phi^{(s-1)}) \leq Q_{\lambda}^{\text{LA}}(\boldsymbol{\beta}^{(s,s-1;k)}, \boldsymbol{\theta}^{(s-1)}, \phi^{(s-1)}).$$

(2) *(Covariance parameter optimization)*

For $l = 1, \dots, d$

$$\theta_l^{(s)} = \arg \min_{\theta_l} Q_{\lambda}^{\text{LA}}(\boldsymbol{\beta}^{(s)}, \boldsymbol{\theta}^{(s,s-1;\theta_l)}, \phi^{(s-1)}).$$

(3) *(Dispersion parameter optimization)*

$$\phi^{(s)} = \arg \min_{\phi} Q_{\lambda}^{\text{LA}}(\boldsymbol{\beta}^{(s)}, \boldsymbol{\theta}^{(s)}, \phi).$$

until convergence.

(i) *Choice of $h_k^{(s)}$.* For $h_k^{(s)}$ we choose the k th diagonal element of the Fisher information of a generalized linear model. Hence, we use the second derivative of the first summand in (7). We set $c_{\min} \leq h_k^{(s)} \leq c_{\max}$ for positive constants c_{\min} and c_{\max} (e.g., $c_{\min} = 10^{-5}$ and $c_{\max} = 10^5$) in order that the algorithm converges (Tseng and Yun 2009).

Algorithm 2 (Approximate) GLMMLasso algorithm

Denote by $\tilde{\mathbf{u}} = \tilde{\mathbf{u}}(\boldsymbol{\beta}^{(s,s-1;k)}, \boldsymbol{\theta}^{(s-1)}, \boldsymbol{\phi}^{(s-1)})$. Replace in Algorithm 1 i)–iii) by

i') Approximate the second derivative

$$\frac{\partial^2}{\partial \beta_k^2} f(\boldsymbol{\beta}^{(s,s-1;\beta_k)}, \boldsymbol{\theta}^{(s-1)}, \boldsymbol{\phi}^{(s-1)} | \tilde{\mathbf{u}}) \Big|_{\beta_k = \beta_k^{(s-1)}}$$

by $h_k^{(s)} > 0$ as described in the section below.

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

$$d_k^{(s)} := \arg \min_d \left\{ f(\boldsymbol{\beta}^{(s,s-1;k)}, \boldsymbol{\theta}^{(s-1)}, \boldsymbol{\phi}^{(s-1)} | \tilde{\mathbf{u}}) + \frac{\partial}{\partial \beta_k} f(\boldsymbol{\beta}^{(s,s-1;\beta_k)}, \boldsymbol{\theta}^{(s-1)}, \boldsymbol{\phi}^{(s-1)} | \tilde{\mathbf{u}}) \Big|_{\beta_k = \beta_k^{(s-1)}} d + \frac{1}{2} d^2 h_k^{(s)} + \lambda \|\boldsymbol{\beta}^{(s,s-1;k)} + d \mathbf{e}_k\|_1 \right\}.$$

iii') Choose a step size $\alpha_k^{(s)} > 0$ and set $\boldsymbol{\beta}^{(s,s-1;k+1)} = \boldsymbol{\beta}^{(s,s-1;k)} + \alpha_k^{(s)} d_k^{(s)} \mathbf{e}_k$ such that

$$Q_\lambda^{\text{LA}}(\boldsymbol{\beta}^{(s,s-1;k+1)}, \boldsymbol{\theta}^{(s-1)}, \boldsymbol{\phi}^{(s-1)} | \tilde{\mathbf{u}}) \leq Q_\lambda^{\text{LA}}(\boldsymbol{\beta}^{(s,s-1;k)}, \boldsymbol{\theta}^{(s-1)}, \boldsymbol{\phi}^{(s-1)} | \tilde{\mathbf{u}}).$$

(ii) *Calculation of $d_k^{(s)}$.* The value $d_k^{(s)}$ is the minimizer of the quadratic approximation of the objective function $Q_\lambda^{\text{LA}}(\cdot)$ and analytically given by Tseng and Yun (2009)

$$d_k^{(s)} = \begin{cases} \text{median} \left(\frac{\lambda - \partial/\partial \beta_k f_{\beta_k}}{h_k^{(s)}}, -\beta_k, \frac{-\lambda - \partial/\partial \beta_k f_{\beta_k}}{h_k^{(s)}} \right) & \text{if } \beta_k \text{ penalized} \\ -\frac{\partial/\partial \beta_k f_{\beta_k}}{h_k^{(s)}} & \text{otherwise,} \end{cases} \quad (9)$$

where $f_{\beta_k} = f(\boldsymbol{\beta}^{(s,s-1;k)}, \boldsymbol{\theta}^{(s-1)}, \boldsymbol{\phi}^{(s-1)})$ in Algorithm 1 and $f_{\beta_k} = f(\boldsymbol{\beta}^{(s,s-1;k)}, \boldsymbol{\theta}^{(s-1)}, \boldsymbol{\phi}^{(s-1)} | \tilde{\mathbf{u}})$ in Algorithm 2.

(iii) *Choice of $\alpha_k^{(s)}$.* The step length $\alpha_k^{(s)}$ is chosen such that the objective function $Q_\lambda^{\text{LA}}(\cdot)$ decreases. We suggest to use the Armijo rule, which is defined for Algorithm 1 as follows (and correspondingly for Algorithm 2 with fixed $\tilde{\mathbf{u}}$):

Armijo rule: Choose $\alpha_k^{\text{init}} > 0$ and let $\alpha_k^{(s)}$ be the largest element of $\{\alpha_k^{\text{init}} \delta^l\}_{l=0,1,2,\dots}$ satisfying

$$Q_\lambda^{\text{LA}}(\boldsymbol{\beta}^{(s,s-1;k)} + \alpha_k^{(s)} d_k^{(s)} \mathbf{e}_k, \boldsymbol{\theta}^{(s-1)}, \boldsymbol{\phi}^{(s-1)}) \leq Q_\lambda^{\text{LA}}(\boldsymbol{\beta}^{(s,s-1;k)}, \boldsymbol{\theta}^{(s-1)}, \boldsymbol{\phi}^{(s-1)}) + \alpha_k^{(s)} \varrho \Delta^k$$

where $\Delta^k := \partial/\partial \beta_k f_{\beta_k} d_k^{(s)} + \gamma (d_k^{(s)})^2 h_k^{(s)} + \lambda \|\boldsymbol{\beta}^{(s,s-1;k)} + d_k^{(s)} \mathbf{e}_k\|_1 - \lambda \|\boldsymbol{\beta}^{(s,s-1;k)}\|_1$.

The choice of the constants comply with the suggestions in Bertsekas (1999), for example, $\alpha_k^{\text{init}} = 1$, $\delta = 0.5$, $\varrho = 0.1$, and $\gamma = 0$.

3.3.2 Active Set Algorithm. If we assume that the true fixed-effect parameter β_0 is sparse in the sense that many elements are zero, we can reduce the computing time remarkably by using an active set algorithm. This is also used by Meier, van de Geer, and Bühlmann (2008), and Friedman, Hastie, and Tibshirani (2010). In particular, we only cycle through all p coordinates every D th iteration, otherwise only through the current active set $S(\hat{\beta}^{(s-1)}) = \{k : \hat{\beta}_k^{(s-1)} \neq 0\}$. Typical values for D are 5 and 10.

An implementation of the algorithm is given in the R package `glmmlmixedlasso` and will be made available on R-Forge (<http://r-forge.R-project.org/>).

4. THE TWO-STAGE GLMMLASSO^{LA} ESTIMATOR(S)

From the soft-thresholding property of the Lasso in linear models (Tibshirani 1996) and in Gaussian linear mixed models (Schelldorfer, Bühlmann, and van de Geer 2011), the fixed-effect estimate $\hat{\beta}$ is biased toward zero. In some GLMMs the estimate of the covariance parameter θ is biased, too. To mitigate these bias problems and the approximation error induced by using the approximate GLMMLasso algorithm, we advocate a two-stage procedure. The first step aims at estimating a candidate set of predictors \hat{S} and can be seen as a variable screening procedure. The purpose of the second step is a more unbiased estimation of the parameters using unpenalized ML estimation based on the selected variables \hat{S} from the first step. The proposed two-stage GLMMLasso algorithm is summarized in Algorithm 3:

Algorithm 3 Two-stage GLMMLasso algorithm

Stage 1: Compute the GLMMLasso^{LA} estimate (8) and the set \hat{S} .

Stage 2: Perform unpenalized ML estimation.

In the next sections, we are going to discuss the specification of the set of variables \hat{S} . We propose two methods from the high-dimensional linear regression framework, and we do not consider the adaptive Lasso (Zou 2006).

4.1 THE GLMMLASSO^{LA}-MLE HYBRID ESTIMATOR

The LARS-OLS hybrid estimator was examined by Efron et al. (2004) and also used by Meinshausen and Bühlmann (2006) and Meier, van de Geer, and Bühlmann (2008). In our context, it becomes a two-stage procedure where the model is refitted including only the covariates with a nonzero fixed-effect coefficient in $\hat{\beta}_{\text{init}}$, where $(\hat{\beta}_{\text{init}}, \hat{\theta}_{\text{init}}, \hat{\phi}_{\text{init}})$ denotes the initial estimate from (8). More specifically, choose $\hat{S} = \hat{S}_{\text{init}} := \{k : |\hat{\beta}_{k,\text{init}}| \neq 0\}$. Then the GLMMLasso^{LA}-MLE hybrid estimator is given by

$$(\hat{\beta}, \hat{\theta}, \hat{\phi})_{\text{hybrid}} := \arg \min_{\beta_{\hat{S}_{\text{init}}}, \theta, \phi} -2 \log L(\beta_{\hat{S}_{\text{init}}}, \theta, \phi), \quad (10)$$

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

4.2 THE THRESHOLDED GLMMLASSO^{LA} ESTIMATOR

The thresholded Lasso with refitting in high-dimensional linear regression models was examined by Zhou (2010) and van de Geer, Bühlmann, and Zhou (2011). We define the set \hat{S}_{thres} to be the set of variables that have initial fixed-effect coefficients larger than some threshold $\lambda_{\text{thres}} > 0$, that is, we choose $\hat{S} = \hat{S}_{\text{thres}} := \{k : |\hat{\beta}_{k,\text{init}}| > \lambda_{\text{thres}}\}$. The thresholded GLMMLasso^{LA} estimator is then defined by

$$(\hat{\beta}, \hat{\theta}, \hat{\phi})_{\text{thres}} := \arg \min_{\beta_{\hat{S}_{\text{thres}}}, \theta, \phi} -2 \log L(\beta_{\hat{S}_{\text{thres}}}, \theta, \phi). \quad (11)$$

The thresholded GLMMLasso^{LA} estimator involves another regularization parameter λ_{thres} , which is determined by minimizing an information criterion presented in the next section.

4.3 SELECTION OF THE REGULARIZATION PARAMETERS

Estimators (8), (10), and (11) require the choice of the regularization parameters λ and λ_{thres} , respectively. We propose to use the Bayesian information criterion (BIC) and the Akaike information criterion (AIC), defined by

$$c_{n,\lambda} = -2 \log L(\hat{\beta}, \hat{\theta}, \hat{\phi}) + a(n) \cdot \hat{d}f_{\lambda} \quad (12)$$

where $a(n) = \log(n)$ for the BIC and $a(n) = 2$ for the AIC. Here, $\hat{d}f_{\lambda} = |\{1 \leq k \leq p : \hat{\beta}_k \neq 0\}| + \dim(\hat{\theta})$ is the sum of the number of nonzero fixed-effect coefficients and the number of covariance parameters. The first summand is motivated by the work of Zou, Hastie, and Tibshirani (2007). The second summand is the approach by Bates (2010), who proposed that in the classical generalized mixed-effects model the degrees of freedom are given by the number of unconstrained optimization parameters. Based on our empirical experience, we suggest for the estimators (8) and (10) the BIC, whereas for (11) we advocate using the AIC (allowing for a larger number of variables) to select λ first and then, sequentially, the BIC to select λ_{thres} . We will compare the performance of the three estimators in the next sections.

5. SIMULATION STUDY

In this section, we assess the performance of the GLMMLasso^{LA} estimators (8), (10), and (11). We compare them with appropriate Lasso, ML, and Penalized Quasi-Likelihood (PQL; Breslow and Clayton 1993) methods.

In the main text, we only present simulation results for the high-dimensional logistic mixed model. Simulation studies for the low-dimensional logistic and the Poisson mixed model are included in the online supplementary materials. At the end of this section, we compare the GLMMLasso^{LA} estimates in a situation where the number of noise variables grows successively.

First of all, let us summarize some general conclusions drawn from real data analysis and the simulation studies:

- (a) The variable screening performance of the GLMMLasso algorithm is not only attractive for the high-dimensional setting, but also for low-dimensional data with a relatively large number of variables (say $p > 20$).
- (b) The GLMMLasso algorithm is numerically as stable as standard R functions like `glmmer` (Bates 2010) or `glmPQL` (Breslow and Clayton 1993; Venables and Ripley 2002) when $p < n$. On the other hand, `glmPath` (Park and Hastie 2007) and `glmnet` (Friedman, Hastie, and Tibshirani 2010) may fail to converge when high-dimensional models are misspecified.
- (c) The main difference between the logistic and the Poisson mixed model is the shrinkage of the covariance parameter estimates of the GLMMLasso^{LA} estimator. These estimates are severely biased in logistic mixed models, in contrast to the Poisson mixed model. Further differences between these two classes are summarized in the online supplementary materials.
- (d) The number of iterations s substantially differs between the classes of GLMMs and the dataset.

5.1 PREVIEW FOR THE LOGISTIC MIXED MODEL

In this section, we confine the discussion to the logistic mixed model because it is viewed as the most challenging model within the class of GLMMs (Molenberghs and Verbeke 2005; Jiang 2007). As an overview, let us sum up the main findings from the simulation study in the logistic mixed model:

- (i) The GLMMLasso^{LA} estimate from (8) of the covariance parameter θ is notably biased. In other words, adding an ℓ_1 -penalty does not only shrink the fixed-effects estimate $\hat{\beta}$, but also the covariance parameter estimate $\hat{\theta}$.
- (ii) In the high-dimensional settings, the GLMMLasso^{LA}-MLE hybrid estimator (10) performs better in terms of parameter estimation accuracy than the thresholded GLMMLasso^{LA} estimator (11).
- (iii) The more the random effects, the more important it is to use the GLMMLasso^{LA} for variable screening (instead of a Lasso ignoring the grouping structure).
- (iv) The number of total iterations s needed is small, often about 15 iterations.

5.2 HIGH-DIMENSIONAL LOGISTIC MIXED MODEL

In all subsequent simulation schemes (including the online supplementary materials), we restrict ourselves to the case where the number of observations per cluster is equal, that is, $n_r = n_C$ for $r = 1, \dots, N$. The covariates are generated from a multivariate normal distribution with mean zero and covariance matrix V with pairwise correlation $V_{kk'} = \rho^{|k-k'|}$ and $\rho = 0.2$. Denote by β_0 the true fixed effects (wherein $(\beta_0)_1$ is the intercept) and by s_0 the true number of nonzero fixed-effect coefficients.

For the logistic mixed models, the intercept and the first covariate have independent random effects with different variance parameters. In particular, $\theta = (\theta_1, \theta_2)$ and covariance

matrix $\Sigma_\theta = \text{diag}(\theta_1^2, \dots, \theta_1^2, \theta_2^2, \dots, \theta_2^2) \in \mathbb{R}^{2N}$, that is, $q = 2N$. We investigate the following two examples in the high-dimensional setting:

H_1 : $N = 40$, $n_C = 10$, $n = 400$, $p = 500$, $\theta_1^2 = \theta_2^2 = 1$ and $s_0 = 5$ with $\beta_0 = (0.1, 1, -1, 1, -1, 0, \dots, 0)^T$.

H_2 : $N = 50$, $n_C = 10$, $n = 500$, $p = 1500$, $\theta_1^2 = \theta_2^2 = 1$ and $s_0 = 5$ with $\beta_0 = (0.1, 1, -1, 1, -1, 0, \dots, 0)^T$.

The fitted models are all correctly specified. Hereafter, we denote by *oracle* the ML estimate of the model that includes only the variables from the true active set. Let *glmmlasso*, *hybrid glmmlasso*, and *thres glmmlasso* be the GLMMLasso^{LA} estimates (8), (10), and (11), respectively. We compare the GLMMLasso^{LA} methods with the standard Lasso for generalized linear models (which ignore the grouping structure). For that purpose, we use the *glm*path algorithm (Park and Hastie 2007) and the BIC as variable selection criterion. Then, let *hybrid glm*path and *thres glm*path be the two-stage procedures based on *glm*path (without random effects).

The results in the form of median and rescaled median absolute deviation (in parentheses) over 100 simulation runs are shown in Table 1. There, $|S(\hat{\beta})|$ denotes the cardinality of the estimated active set and TP is the number of true positives (selected variables that are in the true active set). SE is the squared error of the fixed-effect coefficients, that is, $\text{SE} = \|\hat{\beta} - \beta_0\|_2^2$.

Comparing the cardinality of the active set, we see that *thres glmmlasso* and *thres glm*path have much larger active sets than *glmmlasso* and *glm*path, respectively. This is largely because we employ the AIC in the first and the BIC in the second stage. This is outweighed by the advantage that, on average (not shown), the true effects are predominantly included in *thres glmmlasso*. The active set of *glmmlasso* is slightly smaller than that of *glm*path. And yet, the number of TP is similar as for *glm*path. Hence, we conclude that the existence of random effects does affect the variable selection performance of *glm*path.

Concerning covariance parameter estimation, we read off from the table that $\hat{\theta}_1^2$ and $\hat{\theta}_2^2$ are seriously biased for *glmmlasso*. This motivates the usage of a two-stage procedure. The table suggests that the hybrid and the thresholded procedures have improved estimation accuracy of the random-effects parameters compared to their original counterparts.

Looking at the fixed-effect parameter estimation accuracy, the simulation study reveals that the *glmmlasso* estimates are less biased than the corresponding *glm*path estimates, resulting in lower squared error. And the same holds for *hybrid glmmlasso* and *hybrid glm*path. The fixed-effect parameter estimates of *thres glmmlasso* and *thres glm*path perform inadequately compared to their *hybrid* counterparts. As marked by an asterisk in the table, β_2 is not subject to penalization for the GLMMLasso^{LA} estimator since this variable has a random effect (Schelldorfer, Bühlmann, and van de Geer 2011). Thus, the bias of the estimate is much smaller than for the other fixed-effect coefficients.

To sum up the simulation study, we first conclude that *hybrid glmmlasso* outperforms *thres glmmlasso* in terms of parameter estimation accuracy, with similar performance regarding true positives. Second, *glmmlasso* procedures do outperform *glm*path procedures as variable screening methods.

Of course, *glm*path is fitting a wrong model without random effects.

Table 1. Simulation results (medians) for the logistic mixed models H_1 and H_2 (rescaled median absolute deviations in parentheses). An asterisk (*) means that the corresponding coefficient is not subject to penalization in the GLMMLasso^{LA} estimate

Model	Method	$ S(\hat{\beta}) $	TP	$\hat{\theta}_1^2$	$\hat{\theta}_2^2$	$\hat{\beta}_1^*$	$\hat{\beta}_2^*$	$\hat{\beta}_3$	$\hat{\beta}_4$	$\hat{\beta}_5$	SE	
True		5	5	1	1	0.1	1	-1	1	-1		
H_1	oracle	5 (0)	5 (0)	0.85 (0.4)	0.86 (0.59)	0.07 (0.2)	1.04 (0.25)	-0.99 (0.22)	0.98 (0.18)	-1.01 (0.14)	0.14 (0.088)	
	glmmlasso	6 (1.48)	5 (0)	0.38 (0.24)	0.37 (0.3)	0.06 (0.14)	0.66 (0.16)	-0.3 (0.14)	0.26 (0.14)	-0.34 (0.12)	1.6 (0.42)	
	glmpath	7 (2.22)	5 (0)	- (-)	- (-)	0.04 (0.13)	0.24 (0.12)	-0.21 (0.11)	0.22 (0.1)	-0.28 (0.1)	2.4 (0.52)	
	hybrid glmmlasso	6 (1.48)	5 (0)	0.89 (0.43)	0.87 (0.58)	0.08 (0.19)	1.05 (0.25)	-0.99 (0.23)	1 (0.18)	-1.03 (0.16)	0.44 (0.32)	
	hybrid glmpath	7 (2.22)	5 (0)	0.86 (0.42)	0.87 (0.53)	0.08 (0.2)	1.01 (0.28)	-0.99 (0.24)	0.99 (0.19)	-1.02 (0.16)	0.7 (0.64)	
	thres glmmlasso	10 (3.71)	5 (0)	1.02 (0.7)	1.11 (0.85)	0.1 (0.22)	1.19 (0.29)	-1.09 (0.23)	1.11 (0.2)	-1.13 (0.19)	1.3 (0.77)	
	thres glmpath	10 (2.97)	5 (0)	0.91 (0.49)	0.94 (0.59)	0.09 (0.21)	1.11 (0.27)	-1.07 (0.25)	1.11 (0.19)	-1.1 (0.2)	1.1 (0.73)	
	H_2	oracle	5 (0)	5 (0)	0.89 (0.4)	0.94 (0.53)	0.11 (0.18)	1.02 (0.25)	-0.98 (0.15)	1.02 (0.18)	-1.02 (0.16)	0.13 (0.1)
		glmmlasso	6 (1.48)	5 (0)	0.39 (0.23)	0.41 (0.28)	0.09 (0.13)	0.66 (0.17)	-0.31 (0.1)	0.27 (0.11)	-0.34 (0.09)	1.6 (0.27)
glmpath		6.5 (0.74)	5 (0)	- (-)	- (-)	0.08 (0.11)	0.23 (0.13)	-0.21 (0.08)	0.21 (0.11)	-0.28 (0.08)	2.4 (0.34)	
hybrid glmmlasso		6 (1.48)	5 (0)	0.93 (0.44)	0.96 (0.51)	0.12 (0.19)	1.02 (0.26)	-0.99 (0.15)	1.05 (0.17)	-1.04 (0.16)	0.34 (0.3)	
hybrid glmpath		6.5 (0.74)	5 (0)	0.87 (0.42)	0.94 (0.5)	0.12 (0.18)	1.01 (0.22)	-0.99 (0.15)	1.03 (0.18)	-1.04 (0.17)	0.48 (0.37)	
thres glmmlasso		14 (5.93)	5 (0)	1.3 (0.87)	1.33 (0.79)	0.16 (0.27)	1.26 (0.27)	-1.16 (0.28)	1.2 (0.26)	-1.22 (0.24)	2 (1.7)	
thres glmpath		13.5 (5.19)	5 (0)	0.9 (0.52)	1.03 (0.64)	0.17 (0.24)	1.17 (0.25)	-1.07 (0.19)	1.13 (0.22)	-1.15 (0.21)	1.8 (1.2)	

5.3 LOGISTIC MIXED MODEL WITH A GROWING NUMBER OF NOISE COVARIATES

Here, we assess the performance of *glmmlasso* and *hybrid glmmlasso* when the number of noise variables grows successively. In the low-dimensional setting, we compare them with the ML estimate computed by the R function `glmmer` (denoted by *glmmer*). In addition, let *p-glmmer* be the method that performs variable selection in the following way: Eliminate consecutively (backward selection) all variables with a *p*-value larger than 5% until the final model is attained comprising only significant variables. We compare these four methods in terms of their performance of twice the negative out-of-sample log-likelihood. Let us fix the following random intercept model design: $n = 400$, $N = 40$, $n_C = 10$, $\theta^2 = 1$, $\beta_0 = (0, 1, -1, 1, -1)$. We start with $p = 5$ (no noise variables) and raise the number of variables to $p = 65$. The results over 50 simulation runs are depicted in Figure 1.

The figures show that the negative out-of-sample log-likelihood values for *glmmer* grow polynomial whereas the likelihoods for the other methods remain fairly constant. The increase in *glmmer* stems from the fact that it overfits the model for a growing number of

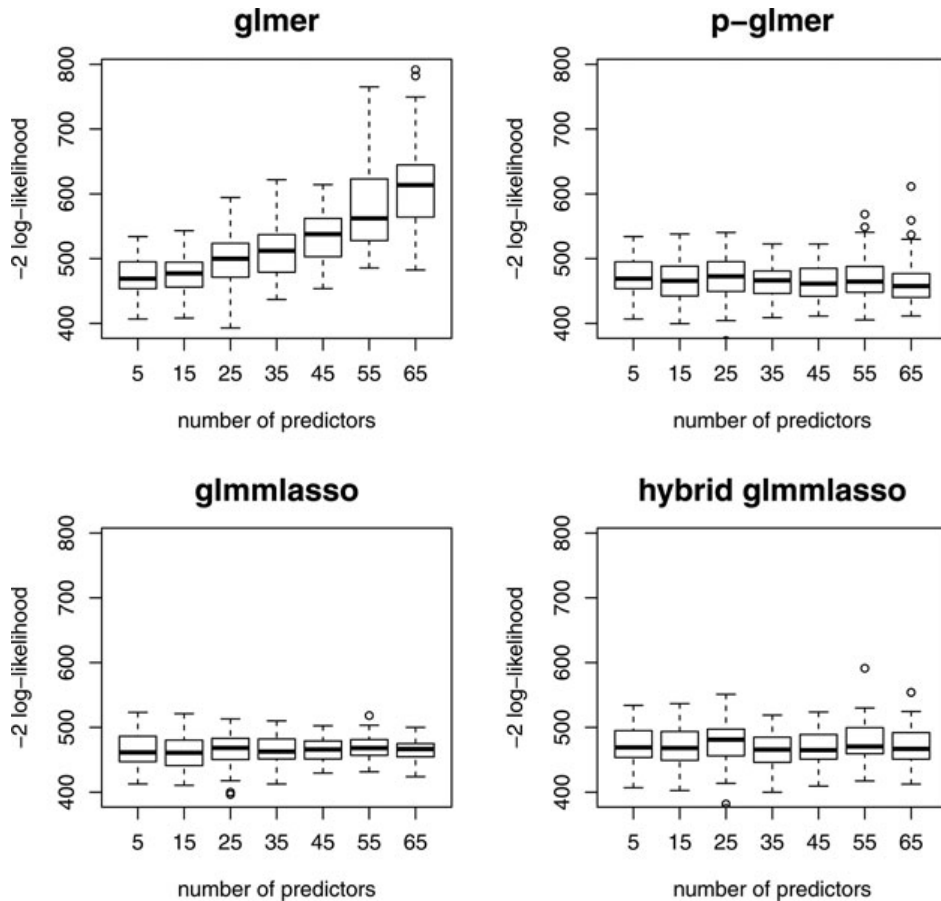


Figure 1. Minus twice out-of-sample log-likelihood for a growing number of covariates. The ML estimate performs badly whereas the GLMMLasso^{LA} estimators remain stable, and they are comparable to the p-glmer in the low-dimensional framework.

covariates. When focusing on the figures in more detail, we read off that the negative log-likelihood of *glmmlasso* increases slightly for larger p whereas the negative log-likelihood of *hybrid glmmlasso* remains stable. The rationale for this small increase in *glmmlasso* is that the more the noise covariates, the larger the optimal λ , and henceforth the larger the shrinkage of the fixed effects. And this leads to the increase of the out-of-sample log-likelihood. *hybrid glmmlasso* (and also *thres glmmlasso*) overcomes this problem and leads to a stable out-of-sample log-likelihood irrespective of p .

5.4 CORRELATED RANDOM EFFECTS

Both from a methodological and an implementational point of view, it is conceptually possible to use correlated random effects. As an illustration we use the logistic mixed model H_1 with correlated random effects (with unstructured covariance matrix) where we use a correlation of $\rho = 0.5$ between the two random effects. The corresponding results are illustrated in Table 2. The results are very similar to the uncorrelated case. However, the bias of the correlation estimate seems to be less severe than the bias of the variance components.

Table 2. Simulation results (medians) for the logistic mixed models H_1 (rescaled median absolute deviations in parentheses). An asterisk (*) means that the corresponding coefficient is not subject to penalization in the GLMMLasso^{LA} estimate

Model	Method	$ S(\hat{\boldsymbol{\beta}}) $	TP	$\hat{\theta}_1^2$	$\hat{\theta}_2^2$	$\hat{\rho}$	$\hat{\beta}_1^*$	$\hat{\beta}_2^*$	$\hat{\beta}_3$	$\hat{\beta}_4$	$\hat{\beta}_5$	SE
True		5	5	1	1	0.5	0.1	1	-1	1	-1	
H_1	oracle	5	5	0.88	0.94	0.53	0.1	0.97	-1.03	1.02	-1.01	0.14
		(0)	(0)	(0.46)	(0.54)	(0.37)	(0.18)	(0.24)	(0.17)	(0.15)	(0.15)	(0.1)
H_1	glmmlasso	6	5	0.41	0.41	0.63	0.07	0.66	-0.33	0.28	-0.34	1.6
		(1.48)	(0)	(0.22)	(0.25)	(0.51)	(0.14)	(0.16)	(0.12)	(0.11)	(0.11)	(0.35)

6. ILLUSTRATION

In this section, we illustrate the proposed GLMMLasso^{LA} estimators for Poisson regression on an extended real dataset with count data.

Data description. We consider the epilepsy data by Thall and Vail (1990) that were also analyzed by Breslow and Clayton (1993). The data were obtained from a randomized clinical trial of 59 patients with epilepsy, comparing a new drug (Trt = 1) with placebo (Trt = 0). The response variable consists of counts of epileptic seizures during the 2 weeks before each of four clinic visits ($V_4 = 1$ for fourth visit, 0 otherwise). Further covariates in the analysis are the logarithm of age (Age), the logarithm of 1/4 the number of baseline seizures (Base), and the interaction of Base and Trt (Base \times Trt). The main question of interest is whether taking the new drug reduces the number of epileptic seizures compared with placebo. To assess the performance of the proposed procedure with high-dimensional data, we add $U(-1, 1)$ distributed noise predictors to get a dataset with $n = 236$, $N = 59$, $n_r = 4$ for $r = 1, \dots, N$, and $p = 4000$. All predictors are standardized to have mean zero and standard deviation one.

Model. Model III in Breslow and Clayton (1993) is a two-level GLMM (Bates 2010), which is an extension of the single-level GLMM introduced in Section 2 for more than one grouping variable. The model consists of two independent random intercept effects. One for subject (Level 1, index r) and one for observation (Level 2, index j). Let θ_{sub}^2 and θ_{obs}^2 be the corresponding variance parameters. Then the linear predictor can be written as

$$\log(\mu_{rj}) = \eta_{rj} = \mathbf{x}_{rj}^T \boldsymbol{\beta} + \theta_{\text{sub}} u_r + \theta_{\text{obs}} u_{rj} \quad r = 1, \dots, 59, \quad j = 1, \dots, 4.$$

Results. The results of the analysis are presented in Table 3. In the first column, we show the estimates for Model III without performing variable selection. There, Intercept, Base, and Trt are significant at the 5% level (indicated by †). If we perform backward selection using the BIC, we end up with a model including Intercept and Base only. And this model coincides with the one selected by *glmmlasso*. *hybrid glmmlasso* overcomes the bias problems of *glmmlasso* and it yields a better model in terms of the BIC. *thres glmmlasso* includes additional noise variables, thereby achieving the smallest BIC score for all models under consideration. Comparing *hybrid glmmlasso* and *thres glmmlasso*, the table suggests that the additional covariates in the latter model reduce the variability while keeping the fixed-effect estimates unaltered.

Table 3. Results for the epilepsy data. Model III is based on six fixed-effect covariates while the other methods are based on $p = 4000$ variables, including 3994 noise covariates. A dagger (\dagger) indicates that the corresponding coefficient is significant at the 5% level. A double dagger (\ddagger) means that five noise variables are selected, but not shown in the table. $S(\hat{\beta}) = \{k : \hat{\beta}_k \neq 0\}$ is the total number of selected variables

	Model III	glmmlasso	hybrid glmmlasso	thres glmmlasso
BIC	527.3	571.8	515.5	480.3
$S(\hat{\beta})$	6	2	2	7 \ddagger
Intercept	1.58 \dagger	1.62	1.58	1.58
Base	0.66 \dagger	$<10^{-4}$	0.74	0.75
Trt	-0.47 \dagger	-	-	-
Base \times Trt	0.36	-	-	-
Age	0.11	-	-	-
V4	-0.04	-	-	-
$\hat{\sigma}_{\text{sub}}^2$	0.21	0.68	0.25	0.28
$\hat{\sigma}_{\text{obs}}^2$	0.13	0.12	0.13	0.04

7. CONCLUDING REMARKS

We address the problem of estimating high-dimensional GLMMs. While low-dimensional GLMMs (Bates 2010) and high-dimensional generalized linear models (van de Geer 2008) have been extensively studied in recent years, little attention has been devoted to high-dimensional GLMMs. We provide an efficient algorithm for the ℓ_1 -penalized ML estimator, called GLMMLasso. It is based on the Laplace approximation, coordinatewise optimization, and a speeding up approximation. The method should be typically used as a screening procedure to estimate a small set of important variables. We propose refitting by ML to get accurate parameter estimates. The second stage is much more important than for linear models, because ℓ_1 -shrinkage can lead to severe bias problems for the estimation of the variance components. Our work is primarily a contribution addressing the numerical challenges of performing high-dimensional variable selection and parameter estimation in nonlinear mixed-effects models involving a nonconvex loss function. An implementation of the algorithm can be found in our R package `glmmlmixedlasso`. It will be made available on R-Forge.

SUPPLEMENTARY MATERIALS

Appendices: Details of the PIRLS algorithm, the comparison of the exact and approximate GLMMLasso algorithms, and additional simulation studies. (`glmmlasso_sm.pdf`)

Dataset: The extended epilepsy dataset used in Section 6. (`epilepsy.txt`)

R-package for GLMMLasso: R-package `glmmlmixedlasso` containing code to perform the GLMMLasso algorithm. (`glmmlmixedlasso-0.1-2.tar.gz`)

ACKNOWLEDGMENTS

The research is supported in part by the Swiss National Science Foundation (grant no. 20PA21-120043/1, ‘‘Forschergruppe FOR 916’’). The authors thank the members of the DFG-SNF Forschergruppe 916 for many

stimulating discussions. In particular, we thank Stephan Dlugosz from the ZEW Mannheim for insisting on studying this particular kind of problem. Moreover, the authors thank the associate editor and two referees for their helpful comments.

[Received September 2011. Revised November 2012.]

REFERENCES

- Bates, D. M. (2010), “lme4: Mixed-Effects Modeling With R,” [online]. Available at <http://lme4.r-forge.r-project.org/book/>. [469,470,474,475]
- (2011a), “Computational Methods for Mixed Models,” [online]. Available at <http://cran.r-project.org/web/packages/lme4/vignettes/Theory.pdf>. [462]
- (2011b), “Linear Mixed Model Implementation in lme4,” [online]. Available at <http://cran.r-project.org/web/packages/lme4/vignettes/Implementation.pdf>. [461,463]
- Bertsekas, D. P. (1999), *Nonlinear Programming*, Belmont, MA: Athena Scientific. [467]
- Breheny, P., and Huang, J. (2011), “Coordinate Descent Algorithms for Nonconvex Penalized Regression, With Applications to Biological Feature Selection,” *Annals of Applied Statistics*, 5, 232–253. [464]
- Breslow, N., and Clayton, D. G. (1993), “Approximate Inference in Generalized Linear Mixed Models,” *Journal of the American Statistical Association*, 88, 9–25. [460,469,470,474]
- Efron, B., Hastie, T., Johnstone, I., and Tibshirani, R. (2004), “Least Angle Regression,” *The Annals of Statistics*, 32, 407–499. [461,468]
- Friedman, J., Hastie, T., and Tibshirani, R. (2010), “Regularization Paths for Generalized Linear Models via Coordinate Descent,” *Journal of Statistical Software*, 33, 1–22. [464,468,470]
- Groll, A., and Tutz, G. (in press), “Variable Selection for Generalized Linear Mixed Models by L_1 -Penalized Estimation,” *Statistics and Computing*. [461]
- Ibrahim, J. G., Zhu, H., Garcia, R. I., and Guo, R. (2010), “Fixed and Random Effects Selection in Mixed Effects Models,” *Biometrics*, 67, 495–503. [461]
- Jiang, J. (2007), *Linear and Generalized Linear Mixed Models and Their Applications*, New York: Springer. [463,470]
- Joe, H. (2008), “Accuracy of Laplace Approximation for Discrete Response Mixed Models,” *Computational Statistics & Data Analysis*, 52, 5066–5074. [463]
- Khalili, A., and Chen, J. (2007), “Variable Selection in Finite Mixture of Regression Models,” *Journal of the American Statistical Association*, 102, 1025–1038. [460]
- Lai, R., Huang, H.-C., and Lee, T. (2012), “Fixed and Random Effects Selection in Nonparametric Additive Mixed Models,” *Electronic Journal of Statistics*, 6, 810–842. [461]
- McCullagh, P., and Nelder, J. A. (1989), *Generalized Linear Models*, London: Chapman & Hall. [460,464]
- McCulloch, C. E., and Searle, S. R. (2001), *Generalized, Linear, and Mixed Models*, New York: Wiley. [460]
- Meier, L., van de Geer, S., and Bühlmann, P. (2008), “The Group Lasso for Logistic Regression,” *Journal of the Royal Statistical Society, Series B*, 70, 53–71. [464,468]
- Meinshausen, N., and Bühlmann, P. (2006), “High-Dimensional Graphs and Variable Selection With the Lasso,” *The Annals of Statistics*, 34, 1436–1462. [468]
- Molenberghs, G., and Verbeke, G. (2005), *Models for Discrete Longitudinal Data*, New York: Springer. [460,463,470]
- Ni, X., Zhang, D., and Zhang, H. H. (2010), “Variable Selection for Semiparametric Mixed Models in Longitudinal Studies,” *Biometrics*, 66, 79–88. [461]
- Pan, W., and Shen, X. (2007), “Penalized Model-Based Clustering With Application to Variable Selection,” *Journal of Machine Learning Research*, 8, 1145–1164. [460]
- Park, M., and Hastie, T. (2007), “ L_1 -Regularization Path Algorithm for Generalized Linear Models,” *Journal of the Royal Statistical Society, Series B*, 69, 659–677. [470,471]

- Schelldorfer, J., Bühlmann, P., and van de Geer, S. (2011), “Estimation for High-Dimensional Linear Mixed-Effects Models Using ℓ_1 -Penalization,” *Scandinavian Journal of Statistics*, 38, 197–214. [461,468,471]
- Skrondal, A., and Rabe-Hesketh, S. (2004), *Generalized Latent Variable Modeling*, Boca Raton, FL: Chapman & Hall/CRC. [463]
- Städler, N., Bühlmann, P., and van de Geer, S. (2010), “ l_1 -Penalization for Mixture Regression Models” (with discussion), *Test*, 19, 209–285. [460]
- Thall, P., and Vail, S. (1990), “Some Covariance Models for Longitudinal Count Data With Overdispersion,” *Biometrics*, 46, 657–671. [474]
- Tibshirani, R. (1996), “Regression Shrinkage and Selection via the Lasso,” *Journal of the Royal Statistical Society, Series B*, 58, 267–288. [460,468]
- Tseng, P., and Yun, S. (2009), “A Coordinate Gradient Descent Method for Nonsmooth Separable Minimization,” *Mathematical Programming, Series B*, 117, 387–423. [461,464,465,467]
- van de Geer, S. (2008), “High-Dimensional Generalized Linear Models and the Lasso,” *The Annals of Statistics*, 36, 614–645. [475]
- van de Geer, S., Bühlmann, P., and Zhou, S. (2011), “The Adaptive and the Thresholded Lasso for Potentially Misspecified Models (and a Lower Bound for the Lasso),” *Electronic Journal of Statistics*, 5, 688–749. [461,469]
- Venables, W. N., and Ripley, B. D. (2002), *Modern Applied Statistics With S*, New York: Springer. [470]
- Witten, D. M., and Tibshirani, R. (2010), “A Framework for Feature Selection in Clustering,” *Journal of the American Statistical Association*, 105, 713–726. [460]
- (2011), “Penalized Classification Using Fisher’s Linear Discriminant,” *Journal of the Royal Statistical Society, Series B*, 73, 753–772. [460]
- Wu, T., and Lange, K. (2008), “Coordinate Descent Algorithms for Lasso Penalized Regression,” *Annals of Applied Statistics*, 2, 224–244. [464]
- Xue, L., Qu, A., and Zhou, J. (2010), “Consistent Model Selection for Marginal Generalized Additive Model for Correlated Data,” *Journal of the American Statistical Association*, 105, 1517–1530. [461]
- Yang, H. (2007), “Variable Selection Procedures for Generalized Linear Mixed Models in Longitudinal Data Analysis,” unpublished Ph.D. dissertation, North Carolina State University. [461]
- Zhou, S. (2010), “Thresholded Lasso for High Dimensional Variable Selection and Statistical Estimation,” *arXiv Preprint arXiv:1002.1583v2*. [461,469]
- Zou, H. (2006), “The Adaptive Lasso and Its Oracle Properties,” *Journal of the American Statistical Association*, 101, 1418–1429. [461,468]
- Zou, H., Hastie, T., and Tibshirani, R. (2007), “On the ‘Degrees of Freedom’ of the Lasso,” *The Annals of Statistics*, 35, 2173–2192. [469]