

1 Hierarchical Selection of Fixed and Random Effects in
2 Generalized Linear Mixed Models

3 Francis K.C. Hui^{*1}, Samuel Müller^{†2}, and A.H. Welsh^{‡1}

4 ¹Mathematical Sciences Institute, The Australian National University,
5 Canberra, Australia

6 ²School of Mathematics and Statistics, University of Sydney, Sydney,
7 Australia

8 **Abstract**

9 In many applications of generalized linear mixed models (GLMMs), there is a hier-
10 archical structure in the effects that needs to be taken into account when performing
11 variable selection. A prime example of this is when fitting mixed models to longitu-
12 dinal data, where it is usual for covariates to be included as only fixed effects or as
13 composite (fixed and random) effects. In this article, we propose the first regularization
14 method that can deal with large numbers of candidate GLMMs while preserving this

*Corresponding author: Francis K.C. Hui, Mathematical Sciences Institute, The Australian National University, 0200, Canberra, ACT, Australia. E: fhui28@gmail.com; P: +61 2 6125 0581.

†E: samuel.mueller@sydney.edu.au

‡E: alan.welsh@anu.edu.au

15 hierarchical structure: CREPE (Composite Random Effects PEnalty) for joint selec-
16 tion in mixed models. CREPE induces sparsity in a hierarchical manner, as the fixed
17 effect for a covariate is shrunk to zero only if the corresponding random effect is or has
18 already been shrunk to zero. In the setting where the number of fixed effects grow at a
19 slower rate than the number of clusters, we show that CREPE is selection consistent for
20 both fixed and random effects, and attains the oracle property. Simulations show that
21 CREPE outperforms some currently available penalized methods for mixed models.

22 **Keywords:** fixed effects, generalized linear mixed models, LASSO, penalized like-
23 lihood, random effects, variable selection

24 1 Introduction

25 Joint selection of fixed and random effects in generalized linear mixed models (GLMMs)
26 presents a challenging problem, especially as regards the question of how to perform selec-
27 tion in a computationally efficient manner while accounting for any hierarchical structure
28 present in the model. Even with a bounded number of covariates, when jointly selecting
29 over fixed and random effects the number of candidate models is considerably larger than in
30 the standard regression context, making methods based on information criteria or the fence
31 (Jiang et al. (2008)) computationally burdensome; see Müller et al. (2013) for a general
32 review of model selection in linear mixed models. One approach to overcoming this compu-
33 tational problem is penalized likelihood methods. While penalized methods for generalized
34 linear models have been extensively studied (dating back to Tibshirani (1996)), their ap-
35 plication to mixed models has only recently been considered, almost exclusively in settings
36 where the number of covariates is bounded, and the selection of fixed and random effects
37 is treated as separate processes. Bondell et al. (2010) and Ibrahim et al. (2011) proposed
38 separate penalties for the fixed and random effects that are summed together. Fan and Li

39 (2012), Peng and Lu (2012), and Lin et al. (2013) all proposed two-stage methods where the
40 fixed and random effects selection are performed independently.

41 When fitting GLMMs to longitudinal data, there is a hierarchical structure in the selection of
42 the effects that is often imposed in practice, namely “we usually only consider time-varying
43 covariates that have been included in the fixed effects.” (Cheng et al. (2010)). It is natural
44 for covariates to be included as either a fixed effect only, or as both fixed and random effects.
45 We refer to the latter as a *composite effect* covariate. As an example, in a longitudinal study
46 monitoring the weights of infants over time (see Section 6), a random slope is included to
47 account for heterogeneity between infants’ changes in weight only if there is a significant
48 overall trend (fixed effect) over time. Another example is in forest management, where
49 random slopes are used to account for between plot variability only if a significant change is
50 observed in the forest’s overall health in response to climate (Hao et al. (2015)). Of course
51 there may be exceptions to this hierarchical structure, a notable one being the case of linear
52 mixed models with centered responses, where a random intercept may be included without a
53 fixed intercept. For most settings however, it is reasonable that covariates should be included
54 as either fixed or composite effects. However, while notions of hierarchical selection have been
55 researched in (generalized) linear models with grouped variables and ordered or polynomial
56 terms, see for instance the group LASSO (Least Absolute Shrinkage and Selection Operator)
57 of Yuan and Lin (2006) and the composite absolute penalty of Zhao et al. (2009), they have
58 not been investigated for GLMMs. This is exemplified in the illustrative examples of Bondell
59 et al. (2010) and Ibrahim et al. (2011), where the respective penalties lead to at least one
60 covariate selected only as a random effect.

61 We propose a penalty called CREPE (Composite Random Effects PEnalty) for hierarchical
62 selection of fixed and random effects in longitudinal GLMMs. CREPE is the first penalty that
63 directly incorporates the notion of covariates being selected as fixed or composite effects. This

64 is done by exploiting the hierarchical structure of the effects, such that a fixed effect coefficient
65 is shrunk to zero only if the corresponding random effect coefficients are, or have already
66 been shrunk to, zero. CREPE also accommodates covariates that are included *a-priori* as
67 fixed effects only. The concept of using a penalty that accounts for the hierarchical structure
68 of the effects has been considered in other contexts, e.g. the fused LASSO (Tibshirani
69 et al. (2005)), finite mixture of regression models (Hui et al. (2015a)), and feature selection
70 in bioinformatics (Garcia et al. (2014)), but has yet to be explored for joint selection in
71 GLMMs. A key part of CREPE’s design involves the use of a group-based penalty for
72 selecting the random effects, specifically, the elements in a row of the eigendecomposition
73 of the random effects covariance matrix (as defined in Section 2) are encouraged to be zero
74 simultaneously.

75 In the setting where the number of fixed effects is allowed to grow at a slower rate than
76 the number of clusters, we show that CREPE satisfies the oracle property of asymptotically
77 identifying the truly non-zero fixed and composite covariates. Regarding computation, we
78 use a Monte-Carlo Expectation Maximization (MCEM, Wei and Tanner (1990)) algorithm to
79 calculate the CREPE estimates, showing how the E-step can be performed straightforwardly
80 for the common cases of Gaussian, Poisson, and Bernoulli responses. Simulation studies show
81 CREPE outperforms some other penalties available for jointly selecting fixed and random
82 effects in GLMMs. We illustrate the application of CREPE to a longitudinal infant study
83 for identifying important baseline and time-varying predictors of infant weights. We provide
84 R code for calculating the CREPE estimates in the Supplementary Material; an R package
85 is planned in future research.

2 Model Selection using CREPE

We focus on the independent cluster model with random intercepts and slopes. Let y_{ij} denote the j^{th} response collected for the i^{th} cluster, where $i = 1, \dots, n$ and $j = 1, \dots, m$. For simplicity, all clusters are assumed to have the same number of measurements, m , where m is bounded and does not grow with n . Conditional on the random effects, the y_{ij} are assumed to be independent responses from the exponential family $f(y_{ij}|\boldsymbol{\beta}, \mathbf{b}_i, \phi)$ with mean μ_{ij} and dispersion parameter ϕ . Given a link function $g(\cdot)$, the mean is modeled as $g(\mu_{ij}) = \boldsymbol{\eta}_{ij} = \mathbf{x}_{ij}^T \boldsymbol{\beta} + \mathbf{z}_{ij}^T \mathbf{b}_i$ for a vector \mathbf{x}_{ij} of predictors corresponding to fixed effects $\boldsymbol{\beta}$, and a vector \mathbf{z}_{ij} of predictors corresponding to random effects \mathbf{b}_i , both containing an intercept term if appropriate. The random effects are assumed to have a multivariate Gaussian distribution, $\mathbf{b}_i \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma})$ where $\boldsymbol{\Sigma} = \boldsymbol{\Gamma} \boldsymbol{\Gamma}^T$ and $\boldsymbol{\Gamma}$ is an unstructured matrix of the same dimension as $\boldsymbol{\Sigma}$, based on the eigendecomposition $\boldsymbol{\Sigma} = \mathbf{Q} \boldsymbol{\Lambda}^{1/2} \boldsymbol{\Lambda}^{1/2} \mathbf{Q}^T = \boldsymbol{\Gamma} \boldsymbol{\Gamma}^T$ such that $\boldsymbol{\Gamma} = \mathbf{Q} \boldsymbol{\Lambda}^{1/2}$, with \mathbf{Q} an orthogonal matrix of normalized eigenvectors and $\boldsymbol{\Lambda}$ a diagonal matrix of eigenvalues.

Lemma 1. *Let $\boldsymbol{\gamma}_k$ be the k^{th} row of $\boldsymbol{\Gamma}$. Then for each k , $\|\boldsymbol{\gamma}_k\| = 0$ implies that $[\boldsymbol{\Sigma}]_{kl} = [\boldsymbol{\Sigma}]_{lk} = 0$ for all l , where $[\boldsymbol{\Sigma}]_{kl}$ refers to element (k, l) of $\boldsymbol{\Sigma}$, and $\|\cdot\|$ denotes the L_2 -norm.*

This result suggests that, rather than penalizing the (diagonal) elements of $\boldsymbol{\Sigma}$ directly, we can employ a group-based penalty on the rows $\boldsymbol{\gamma}_k$, and indeed this is what we pursue. One advantage group-based penalization on the eigendecomposition has is that all the elements of $\boldsymbol{\Gamma}$ can take any number on the real line. This contrasts to the diagonal elements of both $\boldsymbol{\Sigma}$ and its Cholesky decomposition, which are bounded below by zero (see Bondell et al. (2010), Lin et al. (2013), and Pan and Huang (2014) for examples of methods that penalize the diagonal elements of $\boldsymbol{\Sigma}$ or its Cholesky decomposition). By using the eigendecomposition, we can avoid potential boundary issues when performing Taylor expansions (used in the theoretical study of the CREPE estimators in Section 3) and during the actual estimation

110 process.

111 For the independent cluster GLMM, the observed log-likelihood for a GLMM is,

$$\ell(\Psi) = \sum_{i=1}^n \ell_i(\Psi) = \sum_{i=1}^n \log \left(\int \prod_{j=1}^m f(y_{ij} | \beta, \phi, \mathbf{b}_i) f(\mathbf{b}_i | \Gamma) d\mathbf{b}_i \right),$$

112 where $\ell_i(\Psi)$ is the log-likelihood contribution from the i^{th} cluster, and $\Psi = \{\beta, \phi, \text{vec}(\Gamma)\}$.

113 We introduce some notation describing the nature of the covariates in the GLMM. Let α
 114 denote the full set of p covariates in the dataset. We divide this set into mutually exclusive
 115 subsets α_f , which denotes the set of p_f covariates entered into the model as fixed effects
 116 only (e.g., baseline covariates such as gender), and α_c , which denotes the set of p_c covariates
 117 entered into the model as composite effects (e.g., time varying covariates such as time of
 118 visit). We allow p_f to grow at a smaller rate than n (see Condition C6 in Section 3),
 119 while assuming $p_c < m$ is fixed. Subsequently, we can write $\Psi = (\beta, \phi, \gamma_1, \dots, \gamma_{p_c})$ where
 120 $\beta = (\beta_{\alpha_f}, \beta_{\alpha_c})$.

121 The CREPE estimator is defined as the maximizer of the penalized log-likelihood function

$$\ell_{pen}(\Psi) = \ell(\Psi) - n\lambda \sum_{k=1}^p \tilde{w}_k \left(\beta_k^2 + \mathbb{1}_{\{k \in \alpha_c\}} \tilde{v}_k \|\gamma_k\| \right)^{1/2}, \quad (1)$$

122 where $\lambda > 0$ is the tuning parameter and $\mathbb{1}_{\{\cdot\}}$ denotes the indicator function. The adaptive
 123 weights \tilde{w}_k and \tilde{v}_k may depend on a common power parameter $\nu > 0$ (Zou (2006)) and are
 124 required to satisfy some regularity conditions.

125 For $k \in \alpha_f$, CREPE reduces to the adaptive LASSO penalty (Zou (2006)). On the other
 126 hand, for $k \in \alpha_c$, CREPE encourages sparsity in a hierarchical manner so that either both
 127 the fixed and random effects for the covariate are shrunk to zero, or only the random effect
 128 is shrunk to zero. There are two types of sparsity featured in CREPE: group sparsity,
 129 occurring on the rows of the eigendecomposition, $\|\gamma_k\| = 0$, and the “larger” sparsity given

130 by $(\beta_k^2 + \mathbb{1}_{\{k \in \alpha_c\}} \tilde{v}_k \|\gamma_k\|)^{1/2}$. Critically, the group sparsity is nested inside the larger sparsity
 131 event. Thus $\|\gamma_k\| = 0$ must occur either before or simultaneously with $\beta_k = 0$. Then, in
 132 maximizing (1), CREPE allows a covariate $k \in \alpha_c$ to be included as either a fixed effect only,
 133 or as a composite effect.

134 Such a group penalty approach to random effects selection has been considered before by
 135 Ibrahim et al. (2011), and is arguably a better approach than that used by Bondell et al.
 136 (2010) amongst others, which penalizes the diagonal elements of the Cholesky decomposition
 137 of Σ .

138 Fixed intercepts in GLMMs are generally not penalized, although the random intercept (if
 139 included) may be. In such a case, (1) can be altered to $\ell_{pen}(\Psi) = \ell(\Psi) - n\lambda(\tilde{v}_1 \|\gamma_1\|)^{1/2} -$
 140 $n\lambda \sum_{k=2}^p \tilde{w}_k (\beta_k^2 + \mathbb{1}_{\{k \in \alpha_c\}} \tilde{v}_k \|\gamma_k\|)^{1/2}$, where it is assumed the first elements in \mathbf{x}_{ij} and \mathbf{z}_{ij} repre-
 141 sent the fixed and random intercepts respectively.

142 3 Asymptotic Properties

143 We study the large sample properties of the CREPE estimator when p_f grows at a slower
 144 rate than n , while p_c is fixed. Allowing the number of random effects to grow is a more
 145 difficult problem, as it requires both the number of clusters and the cluster size to grow in
 146 order to achieve attractive asymptotic properties (see for instance Fan and Li (2012)), and
 147 (Demidenko (2004)) for an overview of asymptotic theory in mixed models.

148 Let $\Psi_0 = (\beta_0, \phi_0, \gamma_{01}, \dots, \gamma_{0p_c})$, denote the true parameter values, where $\beta_0 = (\beta_{0\alpha_f}, \beta_{0\alpha_c})$
 149 and, let p_{0f} be the number of non-zero elements in $\beta_{0\alpha_f}$. Without loss of generality, we write
 150 $\Psi_0 = (\Psi_{01}, \Psi_{02} = \mathbf{0})$ so Ψ_{01} consists of all the non-zero elements of β_0 , all the vectors γ_{0k}
 151 whose L_2 -norm is positive, and ϕ_0 . Likewise, we write the CREPE estimate as $\hat{\Psi} = (\hat{\Psi}_1, \hat{\Psi}_2)$.
 152 Let $H(\Psi) = -(1/n)\partial^2 \ell(\Psi) / \partial \Psi \partial \Psi^T$ denote the observed Fisher information matrix for the

153 GLMM, and let $\kappa_{\min}\{H(\Psi)\}$ and $\kappa_{\max}\{H(\Psi)\}$ denote its minimum and maximum eigenval-
 154 ues respectively. The following regularity conditions are required here.

155 (C1) For every n , there exists a positive constant c_1 such that $0 < c_1 < \kappa_{\min}\{H(\Psi_0)\} <$
 156 $\kappa_{\max}\{H(\Psi_0)\} < 1/c_1 < \infty$.

157 (C2) For any given $\epsilon > 0$, there exists a $\delta > 0$ with $\|\Psi - \Psi_0\| < \delta$ such that $(1 - \epsilon)c_1 <$
 158 $\kappa_{\min}\{H(\Psi)\} < \kappa_{\max}\{H(\Psi)\} < (1 + \epsilon)/c_1$ for n large enough.

159 (C3) There exists an open subset Ω in the interior of the parameter space of Ψ , containing
 160 Ψ_0 , such that the third derivatives of the log-likelihood $\ell(\Psi)$ exist for every $\Psi \in \Omega$. For
 161 all $\Psi \in \Omega$, there exist integrable functions U_{rst} such that $|\partial^3\ell(\Psi)/\partial\Psi_r\partial\Psi_s\partial\Psi_t| < U_{rst}$,
 162 with $E(U_{rst}^2) < \infty$, where the expectation is with respect to the true model.

163 (C4) $(\min_{l \in \Psi_{01}}\{\beta_{0l}^2\} + \min_{l \in \Psi_{01}}\{\|\gamma_{0l}\|\}) \geq c_2$, where $c_2 > 0$ is a positive constant.

164 (C5) The adaptive weights satisfy $\tilde{w}_k = O_p(1)$ and $\tilde{v}_k = O_p(1)$ for $k \in \Psi_{01}$, and $\tilde{w}_k =$
 165 $O_p\{(n/p_f)^{\nu/2}\}$ and $\tilde{v}_k = O_p\{(n/p_f)^{\nu/2}\}$ for $k \in \Psi_{02}$.

166 (C6) (a) $\lambda\sqrt{np_{0f}} \rightarrow 0$ (b) $\lambda(n/p_f)^{(\nu+3)/4} \rightarrow \infty$, where $\nu > 0$.

167 Condition (C1) ensures the observed Fisher information matrix is well-defined at the true
 168 parameter values for every n , while condition (C2) extends this to a small neighborhood
 169 of Ψ_0 . The two conditions are similar to conditions A4 and A5 in Chen and Chen (2012)
 170 for generalized linear models (GLMs). Condition (C3) is a mild condition to ensure the
 171 log-likelihood function for GLMMs is sufficiently smooth. Since Ψ involves elements of the
 172 eigendecomposition $\mathbf{\Gamma}$ that can take any value on the real line, Ω is guaranteed to not lie
 173 on the boundary space. Condition (C4) places a lower bound on the magnitude of the truly
 174 non-zero coefficients. This may be weakened to permit the truly non-zero effects to tend

175 to zero at a slow rate, although we do not pursue this extension here. Together, conditions
 176 (C2) and (C4) define a rate at which incorrect models are allowed to approach the true
 177 model with increasing n . Condition (C5) is a generalization of condition (C1) in Ibrahim
 178 et al. (2011), requiring that the adaptive weights exhibit different asymptotic behavior for
 179 truly zero and non-zero coefficients. Finally, conditions (C6a) and (C6b) constrain the rate
 180 of growth of the tuning parameter λ , and is similar to conditions in Hui et al. (2015b) for
 181 adaptive LASSO GLMs. Together, they restrict the number of fixed effects to grow subject
 182 to $(p_f/n)^{(\nu+3)/4} \sqrt{np_{0f}} \rightarrow 0$. This is an advance on Ibrahim et al. (2011) and Lin et al. (2013),
 183 amongst others, who proved oracle properties assuming fixed p .

184 We first establish a result regarding the consistency properties of the CREPE estimator.

185 **Theorem 1.** *If (C1)-(C6) are satisfied and $\nu \geq 1$, then there exists a local maximizer $\hat{\Psi}$ of*
 186 *the penalized log-likelihood function in (1) that satisfies*

187 (a) *Estimation consistency: $\|\hat{\Psi} - \Psi_0\| = O_p(\sqrt{p_f/n})$.*

188 (b) *Selection consistency: $P(\hat{\Psi}_2 = \mathbf{0}) \rightarrow 1$.*

189 With probability tending to one then, CREPE asymptotically correctly determines whether
 190 each covariate is a fixed or a composite effect.

191 Let $\mathcal{I}(\Psi_0) = E(-\partial^2 \ell(\Psi) / \partial \Psi \partial \Psi^T) |_{\Psi_0}$ be the expected Fisher information matrix evaluated
 192 at the true parameter point.

193 **Theorem 2.** *For a fixed integer q , let \mathbf{B}_n be a $q \times \dim(\Psi_{01})$ matrix such that $\mathbf{B}_n \mathbf{B}_n^T \rightarrow \mathbf{G}$*
 194 *for some non-negative, symmetric $q \times q$ matrix \mathbf{G} . If (C1)-(C6) are satisfied and $\nu \geq 1$,*
 195 *then the local maximizer $\hat{\Psi}$ in Theorem 1 satisfies*

$$\sqrt{n} \mathbf{B}_n \mathcal{I}^{-1/2}(\Psi_{01})(\hat{\Psi}_1 - \Psi_{01}) \xrightarrow{d} \mathcal{N}(\mathbf{0}, \mathbf{G}),$$

196 where $\mathcal{I}(\Psi_{01})$ is the block of the expected Fisher information matrix involving only the truly
 197 non-zero parameters Ψ_{01} .

198 Theorems 1 and 2 establish that the CREPE estimator attains the oracle property in
 199 GLMMs. The proofs of the theorems are provided in the Supplementary Material, following
 200 a similar outline to that of Fan and Peng (2004).

201 4 Estimation

202 We use the Monte-Carlo EM (MCEM, Wei and Tanner (1990)) algorithm combined with the
 203 local quadratic approximation (Fan and Li (2001)) for calculating the CREPE estimators.
 204 We focus on the common cases of Gaussian, Poisson, and Bernoulli mixed models, showing
 205 that updates of the parameters in these cases can be obtained straightforwardly. Let

$$\begin{aligned} \ell_{pen,c}(\Psi, \mathbf{b}) &= \sum_{i=1}^n \left(\sum_{j=1}^m \log\{f(y_{ij}|\boldsymbol{\beta}, \phi, \mathbf{b}_i)\} - \frac{1}{2} \log\{\det(\mathbf{\Gamma}\mathbf{\Gamma}^T)\} - \frac{1}{2} \mathbf{b}_i^T (\mathbf{\Gamma}\mathbf{\Gamma}^T)^{-1} \mathbf{b}_i \right) \\ &\quad - n\lambda \sum_{k=1}^p \rho(\beta_k, \gamma_k) \\ &= \sum_{i=1}^n \ell_{c,i}(\Psi, \mathbf{b}_i) - n\lambda \sum_{k=1}^p \rho(\beta_k, \gamma_k) \end{aligned}$$

206 where $\rho(\beta_k, \gamma_k) = \tilde{w}_k(\beta_k^2 + \mathbb{1}_{\{k \in \alpha_c\}} \tilde{v}_k \|\boldsymbol{\gamma}_k\|)^{1/2}$. Suppose at iteration t , we have estimates $\hat{\Psi}^{(t)}$.
 207 The MCEM algorithm iterates between the following steps: the E-step, which calculates the
 208 expectation of $\ell_{pen,c}(\Psi, \mathbf{b})$ with respect to the conditional posterior distribution $f(\mathbf{b}_i|\mathbf{y}, \hat{\Psi}^{(t)})$,
 209 better known as the Q-function, and the M-step, which maximizes the Q-function to obtain
 210 updated estimates $\hat{\Psi}^{(t+1)}$. For non-Gaussian responses where the posterior distribution does

211 not possess a closed form, we perform the E-step using Monte-Carlo integration,

$$\begin{aligned} \mathbb{E}_{\mathbf{b}_i|\hat{\Psi}^{(t)}} \{ \ell_{c,i}(\Psi, \mathbf{b}_i) \} &= \int \ell_{c,i}(\Psi, \mathbf{b}_i) \times \frac{\prod_{j=1}^m f(y_{ij}|\hat{\beta}^{(t)}, \hat{\phi}^{(t)}, \mathbf{b}_i) f(\mathbf{b}_i|\hat{\Gamma}^{(t)})}{\exp\{\ell_i(\hat{\Psi}^{(t)})\}} d\mathbf{b}_i \\ &\approx \exp\{\ell_i(\hat{\Psi}^{(t)})\}^{-1} \frac{1}{D} \sum_{d=1}^D \ell_{c,i}(\Psi, \mathbf{b}_i^d) \prod_{j=1}^m f(y_{ij}|\hat{\beta}^{(t)}, \hat{\phi}^{(t)}, \mathbf{b}_i^d), \end{aligned} \quad (2)$$

212 where \mathbf{b}_i^d is simulated from $f(\mathbf{b}_i|\hat{\Gamma}^{(t)})$, the quantity $\exp\{\ell_i(\hat{\Psi}^{(t)})\}$ is approximated as
 213 $D^{-1} \sum_{d=1}^D \prod_{j=1}^m f(y_{ij}|\hat{\beta}^{(t)}, \hat{\phi}^{(t)}, \mathbf{b}_i^d)$, and D is the number of Monte-Carlo samples. In the simula-
 214 tions in Section 5, we used $D = 2,000$.

215 To avoid non-differentiability at the origin, we approximate the CREPE penalty by a local
 216 quadratic approximation (LQA). At iteration t , set element k of $\hat{\Psi}^{(t+1)}$ to zero if the corre-
 217 sponding element in $\hat{\Psi}^{(t)}$ is equal to or very close to zero, e.g., absolute value within 10^{-3} .
 218 Otherwise, approximate the CREPE penalty as

$$\rho(\beta_k, \gamma_k) = \rho(\hat{\beta}_k^{(t)}, \hat{\gamma}_k^{(t)}) + M_k^{(t)}(\beta_k^2 - (\hat{\beta}_k^{(t)})^2) + \mathbb{1}_{\{k \in \alpha_c\}} M_k^{(t)} \frac{\tilde{v}_k}{2\|\hat{\gamma}_k^{(t)}\|} (\gamma_k^T \gamma_k - (\hat{\gamma}_k^{(t)})^T \hat{\gamma}_k^{(t)}),$$

219 where $M_k^{(t)} = (\tilde{w}_k/2) \left((\hat{\beta}_k^{(t)})^2 + \mathbb{1}_{\{k \in \alpha_c\}} \tilde{v}_k \|\hat{\gamma}_k^{(t)}\| \right)^{-1/2}$. Combining these results, the M-step
 220 consists of maximizing the penalized Q-function,

$$Q_{pen}(\Psi|\hat{\Psi}^{(t)}) = \mathbb{E}_{\mathbf{b}_i|\hat{\Psi}^{(t)}} \{ \ell_{c,i}(\Psi, \mathbf{b}_i) \} - n\lambda \sum_{k=1}^p \left(M_k^{(t)} \beta_k^2 + \mathbb{1}_{\{k \in \alpha_c\}} M_k^{(t)} \frac{\tilde{v}_k}{2\|\hat{\gamma}_k^{(t)}\|} \gamma_k^T \gamma_k \right).$$

221 We now focus on the three special cases of Gaussian, Poisson, and Bernoulli responses.
 222 Gaussian responses: For the linear mixed model where $f(y_{ij}|\beta, \phi, \mathbf{b}_i) = \mathcal{N}(\eta_{ij}, \sigma^2)$, a closed
 223 form for the posterior distribution of \mathbf{b}_i can be obtained. Let $\mathbf{y}_i = (y_{i1}, \dots, y_{im})$, $\mathbf{X}_i =$
 224 $(\mathbf{x}_{i1} \dots \mathbf{x}_{im})^T$ and $\mathbf{Z}_i = (\mathbf{z}_{i1} \dots \mathbf{z}_{im})^T$. It is straightforward to show that $f(\mathbf{b}_i|\mathbf{y}, \hat{\Psi}) =$
 225 $\mathcal{N}(\hat{\mathbf{a}}_i, \hat{\mathbf{A}}_i)$, where $\hat{\mathbf{A}}_i = \left((\hat{\Gamma}\hat{\Gamma}^T)^{-1} + \hat{\sigma}^{-2} \mathbf{Z}_i^T \mathbf{Z}_i \right)^{-1}$ and $\hat{\mathbf{a}}_i = \hat{\sigma}^{-2} \hat{\mathbf{A}}_i \mathbf{Z}_i^T (\mathbf{y}_i - \mathbf{X}_i \hat{\beta})$. In turn,
 226 we can derive a closed form for the penalized Q-function by using this result and the fact

227 that

$$E_{\mathbf{b}_i|\hat{\Psi}^{(t)}}(\mathbf{b}_i^T(\mathbf{\Gamma}\mathbf{\Gamma}^T)^{-1}\mathbf{b}_i) = \hat{\mathbf{a}}_i^T(\mathbf{\Gamma}\mathbf{\Gamma}^T)^{-1}\hat{\mathbf{a}}_i + \text{tr}\{(\mathbf{\Gamma}\mathbf{\Gamma}^T)^{-1}\hat{\mathbf{A}}_i\}, \quad (3)$$

228 an identity that does not require the normality assumption on \mathbf{b}_i . Closed form updates for
 229 $\boldsymbol{\beta}$ and σ^2 may then be obtained, while a Quasi-Newton method, for instance, can be used
 230 to update the rows of $\mathbf{\Gamma}$.

231 Poisson responses: Using the log link, we have

232 $\sum_{i=1}^n \sum_{j=1}^m \log\{f(y_{ij}|\boldsymbol{\beta}, \mathbf{b}_i)\} = \sum_{i=1}^n \sum_{j=1}^m \{y_{ij}(\mathbf{x}_{ij}^T\boldsymbol{\beta} + \mathbf{z}_{ij}^T\mathbf{b}_i) - \exp(\mathbf{x}_{ij}^T\boldsymbol{\beta}) \exp(\mathbf{z}_{ij}^T\mathbf{b}_i)\}$. From this, it
 233 is straightforward to see that for the penalized Q-function, we only require Monte-Carlo
 234 estimates of the posterior mean $E_{\mathbf{b}_i|\hat{\Psi}^{(t)}}(\mathbf{b}_i)$, the moment generating function $E_{\mathbf{b}_i}\{\exp(\mathbf{z}_{ij}^T\mathbf{b}_i)\}$,
 235 along with the posterior covariance matrix for use in (3). Since none of these is a function
 236 of the parameters that need updating, the M-step can be performed relatively quickly.

237 Bernoulli responses: Using the logit link, we have

238 $\sum_{i=1}^n \sum_{j=1}^m \log\{f(y_{ij}|\boldsymbol{\beta}, \mathbf{b}_i)\} = \sum_{i=1}^n \sum_{j=1}^m [y_{ij}(\mathbf{x}_{ij}^T\boldsymbol{\beta} + \mathbf{z}_{ij}^T\mathbf{b}_i) - \log\{1 + \exp(\mathbf{x}_{ij}^T\boldsymbol{\beta} + \mathbf{z}_{ij}^T\mathbf{b}_i)\}]$. Applying
 239 the MCEM algorithm directly is challenging because the second term is non-linear in $\boldsymbol{\beta}$. To
 240 overcome this, we use the fact that the variance of the Bernoulli distribution is bounded above
 241 by 1/2. We can therefore minorize the above expression by a partial quadratic expansion
 242 about $\boldsymbol{\beta} = \hat{\boldsymbol{\beta}}^{(t)}$,

$$\begin{aligned} \sum_{i=1}^n \sum_{j=1}^m \log\{f(y_{ij}|\boldsymbol{\beta}, \mathbf{b}_i)\} &\geq \sum_{i=1}^n \sum_{j=1}^m \log\{f(y_{ij}|\hat{\boldsymbol{\beta}}^{(t)}, \mathbf{b}_i)\} + \sum_{i=1}^n \sum_{j=1}^m (y_{ij} - \mu_{ij}^{(t)})\mathbf{x}_{ij}^T(\boldsymbol{\beta} - \hat{\boldsymbol{\beta}}^{(t)}) \\ &\quad - \frac{1}{4} \sum_{i=1}^n \sum_{j=1}^m (\boldsymbol{\beta} - \hat{\boldsymbol{\beta}}^{(t)})^T \mathbf{x}_{ij} \mathbf{x}_{ij}^T (\boldsymbol{\beta} - \hat{\boldsymbol{\beta}}^{(t)}), \end{aligned} \quad (4)$$

243 where $\eta_{ij}^{(t)} = \mathbf{x}_{ij}^T\hat{\boldsymbol{\beta}}^{(t)} + \mathbf{z}_i^T\mathbf{b}_i$ and $\mu_{ij}^{(t)} = \exp(\eta_{ij}^{(t)})/\{1 + \exp(\eta_{ij}^{(t)})\}$ (see Hunter and Li (2005)
 244 for details on the notion of minorizing functions). Since this inequality remains true when
 245 we apply expectations to both sides, it means that we can use (4) to construct a minorizer

246 of $Q_{pen}(\Psi|\hat{\Psi}^{(t)})$, and therefore maximize the minorizer instead. This is known as a (Monte-
 247 Carlo) minorization-maximization algorithm, as detailed in Hunter and Li (2005). Import-
 248 tantly, it is clear that this minorizer requires only Monte-Carlo estimates of $E_{\mathbf{b}_i|\hat{\Psi}^{(t)}}(\mathbf{b}_i)$, the
 249 expected fitted probability $E_{\mathbf{b}_i}(\mu_{ij}^{(t)})$, along with the posterior covariance matrix for use in
 250 (3). As none of these is a function of the parameters that need updating, the maximization
 251 can be performed straightforwardly.

252 5 Simulation Study

253 An empirical study was conducted to compare the performance of CREPE with some other
 254 proposed penalties for variable selection in GLMMs. We focus on the cases of Gaussian,
 255 Poisson and Bernoulli responses. For brevity, only the results for Gaussian and Bernoulli
 256 mixed models are presented; the results for Poisson GLMMs are similar and are provided
 257 in the Supplementary Material. For CREPE, we chose the adaptive weights as follows. Let
 258 $\tilde{\boldsymbol{\beta}} = (\tilde{\boldsymbol{\beta}}_f, \tilde{\boldsymbol{\beta}}_c)$ and $\tilde{\boldsymbol{\Sigma}}$ denote the maximum likelihood estimators of the fixed effects coefficients
 259 and random effects covariance matrix, based on fitting a saturated GLMM using the `lme4`
 260 package (Bates et al. (2014)). Then we set $\tilde{w}_k = |\tilde{\beta}_k|^{-2}$ and $\tilde{v}_k = [\tilde{\boldsymbol{\Sigma}}]_{kk}^{-2}$, where $[\tilde{\boldsymbol{\Sigma}}]_{kk}$ denotes
 261 the k^{th} diagonal element of $\tilde{\boldsymbol{\Sigma}}$. The saturated GLMM fit was also used to obtain starting
 262 values for the CREPE estimator. It is worth pointing out that the current version of `lme4`
 263 (version 1.1-10 at the time of writing) does not permit fitting mixed models when the number
 264 of random effects exceeds cluster size, $p_c > m$. Instead, we used an older version (version
 265 1.0-6) that did permit such saturated models to be fitted.

266 In all three settings, we used a BIC-type criterion to select the tuning parameter for CREPE,
 267 $\text{BIC}_\lambda = -2\ell(\hat{\Psi}) + \log(n) \dim(\hat{\Psi})$, where $\dim(\hat{\Psi})$ denotes the number of *non-zero* estimated
 268 parameters in $\hat{\Psi}$. The model complexity penalty used in the BIC is based on the log of the

269 number of clusters, n . More generally, our use of a BIC-type criterion for tuning parameter
270 selection is comparable to what has been advocated in Bondell et al. (2010) and Lin et al.
271 (2013), amongst others. We did however also consider the use of an AIC-type criterion,
272 where $\log(n)$ was replaced by 2 as the model complexity penalty, with results (not shown)
273 indicating that it tended to overfit both the fixed and random effects.

274 For each combination of n (number of clusters) and m (cluster size) considered, we generated
275 200 datasets. We assessed performance in terms of both model selection and model accuracy.
276 For the former, we considered the mean number of false positives (truly zero coefficients
277 not shrunk to zero, indicative of overfitting) and false negatives (truly non-zero coefficients
278 shrunk to zero, indicative of underfitting) for the fixed effects, and the percentage of datasets
279 with correctly chosen random effects. We also recorded the percentage of datasets where the
280 method produced non-hierarchical shrinkage, where one or more covariates end up being
281 selected as a random effect only. As discussed below (1), such non-hierarchical shrinkage is
282 not permitted by the design of the CREPE penalty. In the Supplementary Material, we also
283 present the percentage of datasets where the method obtained the correct model.

284 To assess model accuracy, we computed two measures for each method: the Kullback-Leibler
285 distance between the true and fitted models, and the model error defined as the squared
286 Euclidean norm between the estimated and true parameters. We subsequently computed a
287 median relative Kullback-Leibler distance and the median relative model error, the median
288 of the ratios of the Kullback-Leibler distance (or model error) between the CREPE estimator
289 and the alternative method. Relative Kullback-Leibler distances and model errors less than
290 one were indicative of CREPE having better model accuracy. Similar measures of model
291 accuracy were used in Bondell et al. (2010) and Lin et al. (2013), among many others.
292 Because the results for both measures were similar, we only present the relative Kullback-
293 Leibler distance results in main text, and present the results for relative model errors in the

294 Supplementary Material.

295 5.1 Normal Responses

296 We adapted the simulation design in Bondell et al. (2010), but allowed the number of fixed
 297 effects to grow with n . In detail, datasets were simulated from a linear mixed model with
 298 the number of predictors growing at rate $p = \lceil 7n^{1/4} \rceil$ where $\lceil \cdot \rceil$ is the ceiling function.
 299 Covariates \mathbf{x}_{ij} were constructed by setting the first element to one for a fixed intercept, and
 300 generating the remaining elements from a multivariate Gaussian distribution with mean zero
 301 and covariance matrix $\text{Cov}(x_{ijr}, x_{ijs}) = 0.5^{|r-s|}$. The covariates for the random effects \mathbf{z}_{ij}
 302 were taken as the first eight covariates of \mathbf{x}_{ij} , so $p_c = 8$ and $p_f = p - p_c$ grows at the same
 303 rate as p . For the true model, the first eight elements of $\boldsymbol{\beta}_0$ were set to $(-1, 3, 1.5, 0, 0, 2, 1, 0)$.
 304 Then every third term was set to alternating values of ± 1 . The true 8×8 covariance matrix
 305 for the random effects, $\boldsymbol{\Sigma}_0$, consisted of two non-zero blocks: I) a 2×2 matrix with diagonal
 306 entries 9 and 4, and off-diagonal entries of 4.8, occupying rows/columns 1 and 2 of $\boldsymbol{\Sigma}_0$, II) a
 307 2×2 diagonal matrix with entries 2, occupying rows/columns 6 and 7 of $\boldsymbol{\Sigma}_0$. This resulted
 308 in four informative composite effect covariates. Responses y_{ij} were then generated from a
 309 Gaussian distribution with variance $\sigma_0^2 = 1$. We considered combinations of $n = 30, 60$
 310 clusters (corresponding to $p = 17$ and 20 respectively) and cluster sizes of $m = 5, 10, 20$.
 311 Three penalized estimators were compared: (1) CREPE with $\nu = 2$ in the adaptive weights
 312 for CREPE, (2) the M-ALASSO penalty of Bondell et al. (2010), and (3) the ALASSO
 313 penalty of Lin et al. (2013). To the best of our knowledge, these three procedures are
 314 currently the only penalties publicly available in R for selecting both fixed and random effects,
 315 and we found no additional methods. Since all procedures perform joint selection of fixed
 316 and random effects, we took the model error as $\text{ME} = \|\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}_0\|^2 + \|\text{vech}(\hat{\boldsymbol{\Sigma}}) - \text{vech}(\boldsymbol{\Sigma}_0)\|^2$.
 317 Overall, CREPE performed the best in selecting both fixed and random effects, as well as

318 in model accuracy (Table 1 and Supplementary Material Table 1). M-ALASSO tended to
 319 choose a smaller number of fixed effects compared to CREPE, as reflected in the lower number
 320 of false positives but higher number of false negatives, while ALASSO performed worst as it
 321 severely overfitted the fixed effects. For random effects, M-ALASSO performed slightly better
 322 than CREPE although differences between the two were minor at the larger cluster sizes. For
 323 all settings, CREPE performed best in terms of selecting the correct model (Supplementary
 324 Material Table 1). ALASSO tended to underfit the random effects and shrink rows/columns
 325 6 and 7 of the covariance matrix to zero. This underfitting of the random effects by ALASSO
 326 may be a result of the BIC used for the selecting the tuning parameter, which involves a large
 327 model complexity penalty $\log(mn)$ (following the recommendation in Lin et al., 2013). The
 328 median relative Kullback-Leibler distance was less than one in all but one case, indicating
 329 that CREPE has better model accuracy compared to the two alternative methods.

330 Both M-ALASSO and ALASSO presented cases of non-hierarchical shrinkage, particularly
 331 on element 7 in \mathbf{x}_{ij} (and equivalently \mathbf{z}_{ij}) where the fixed effect was shrunk to zero while the
 332 corresponding random effect remained in the final model. Not surprisingly, the percentage
 333 of datasets where non-hierarchical shrinkage occurred decreased with increasing cluster size
 334 m .

335 5.2 Bernoulli Responses

336 We generated datasets from a Bernoulli GLMM using the same rate of growth of p (and thus
 337 p_f) as in Section 5.1. Covariates \mathbf{x}_{ij} and \mathbf{z}_{ij} were constructed in the same manner as in the
 338 Gaussian response case, \mathbf{z}_{ij} being taken as the first eight covariates of \mathbf{x}_{ij} such that $p_c = 8$.
 339 The elements of $\boldsymbol{\beta}_0$ were the same as in Setting 1, while the true 8×8 covariance matrix $\boldsymbol{\Sigma}_0$
 340 was set to a diagonal matrix with the entries $(1, 1, 0, 0, 0, 1, 0, 0)$. Responses y_{ij} were then
 341 generated from a Bernoulli distribution with logit link. For CREPE, we used $\nu = 2$ for the

Table 1: Simulation results for linear mixed models. Performance was assessed the mean number false positives (FP) and false negatives (FN) for the fixed effects, the percentage of datasets with correctly chosen random effects components (%RE), percentage of datasets where there was non-hierarchical shrinkage (%S), and median relative Kullback-Leibler distance (RKL). Since %S was equal to zero for CREPE, this column is omitted from the table.

n	m	CREPE			M-ALASSO					ALASSO				
		FP	FN	%RE	FP	FN	%RE	%S	RKL	FP	FN	%RE	%S	RKL
30	5	0.52	0.19	38	0.23	1.02	47	78	0.92	3.21	0.62	4	94	0.83
	10	0.05	0.06	86	0.03	0.28	90	29	0.90	2.45	0.53	50	50	0.78
	20	0.06	0.02	95	0.01	0.24	96	24	0.50	4.46	0.42	41	35	0.39
60	5	0.32	0.03	42	0.05	0.28	63	47	0.82	1.09	0.34	38	76	1.01
	10	0	0.02	93	0	0.10	94	14	0.64	1.44	0.39	72	40	0.95
	20	0.01	0	97	0.01	0.07	96	9	0.49	3.37	0.31	56	39	0.63

342 adaptive LASSO weights. We considered combinations of $n = 50, 100$ clusters, corresponding
343 to $p = 19$ and 23 respectively, and cluster sizes of $m = 10, 20$. We had intended to perform
344 simulations at $m = 5$ also, as we did with Gaussian and Poisson responses, but found that we
345 were unable to obtain suitable adaptive weights for CREPE based on a saturated GLMM fit.
346 This was not surprising given the small cluster size $m = 5$ and relative lack of information in
347 Bernoulli responses. While other methods of obtaining adaptive weights are possible, they
348 are outside the scope of this work (see also our discussion in Section 7).

349 To our knowledge, no R packages are currently available for performing joint selection in
350 mixed models with non-normal responses. For comparison with CREPE then, we considered
351 the `glmmLasso` package (Groll and Tutz (2014)), which performs fixed effects selection only
352 in GLMMs using the unweighted LASSO penalty. With this method, we considered two
353 possibilities: the random effects component was known and only elements 1, 2, and 6 of \mathbf{z}_{ij}
354 were included; the random effects was unknown and all eight elements of \mathbf{z}_{ij} were included.
355 Our fitting models of such models via `glmmLasso` is unconventional in allowing fixed effects

356 to be penalized when the corresponding random effects (by definition of the program) cannot
357 be penalized. We see this less as an argument against `glmLasso` and more one in favour of
358 using CREPE as a penalty.

359 Because `glmLasso` only performs selection of the fixed effects here, the model error is based
360 only on the fixed effects, $ME = \|\hat{\beta} - \beta_0\|^2$. This avoids confounding the results with whether
361 the true and saturated random effects structure was used for `glmLasso`. We considered
362 several ways of implementing the package, and we present results based on the method
363 which worked best, namely constructing a solution path from the smallest to the largest
364 value of the tuning parameter.

365 CREPE performed better than both versions of `glmLasso` at selecting the fixed effects,
366 except at $n = 50$ and $m = 10$ where it had a slight tendency to underfit the fixed effects
367 (Table 2 and Supplementary Material Table 3). This underfitting may explain why the
368 relative Kullback-Leibler distance for both versions of `glmLasso` was greater than one for
369 this setting. In all other settings, CREPE had better model accuracy as reflected in the
370 relative Kullback-Leibler distance (and model errors in Supplementary Material Table 2).
371 At $n = 50$, both versions of `glmLasso` tended to overfit the fixed effects, a result that may
372 be partly attributed to the lack of adaptive weights. Regarding random effects selection,
373 even at $n = 100$ and $m = 20$, CREPE was only able to correctly pick the true random effects
374 structure half the time, with a tendency to overfit and fail to shrink rows/column 3 of the
375 estimated \mathbf{D} to zero (note this covariate has a corresponding non-zero fixed effect).

376 When the true random effects structure was known, `glmLasso` presented no cases of non-
377 hierarchical shrinkage (%S). By contrast, when a saturated structure was assumed for the
378 random effects, strong evidence of non-hierarchical shrinkage was observed for `glmLasso`,
379 as it shrank one or more of the fixed effects for covariates 4, 5, and 8 to zero while leaving
380 the corresponding random effects in the model. This was not surprising as our application

381 of `glmLasso` allows fixed effects to be penalized in a situation where the program (by
 382 definition) cannot penalize the corresponding random effects.

Table 2: Simulation results for Bernoulli GLMMs. Performance was assessed based on the mean number false positives (FP) and false negatives (FN) for the fixed effects, the percentage of datasets with correctly chosen random effects components (%RE, for CREPE only), the percentage of datasets where there was non-hierarchical shrinkage (%S), and median relative Kullback-Leibler distance (RKL). Since %S was equal to zero for CREPE, the column is omitted from the table.

n	m	CREPE			<code>glmLasso_{true}</code>				<code>glmLasso_{sat}</code>			
		FP	FN	%RE	FP	FN	%S	RKL	FP	FN	%S	RKL
50	10	0.68	0.71	17	1.44	0.06	0	1.18	1.55	0.05	96	1.18
	20	0.13	0.01	31	2.54	0	0	0.74	3.55	0	87	0.70
100	10	0.15	0.02	11	0.57	0	0	0.85	0.78	0	100	0.82
	20	0.04	0	51	0.34	0	0	0.55	0.47	0	100	0.56

383 6 Application to Yale Infant Study

384 To illustrate the application of CREPE, we analyzed the Yale infant growth study of Wasser-
 385 man and Leventhal (1993), which aimed to identify, among other things, whether cocaine
 386 exposure during pregnancy affects weight gain in children. The dataset was also used in
 387 Ibrahim et al. (2011). A total of $n = 298$ infants were recruited for the study, and their
 388 weight (in pounds) monitored over the study period. Seven predictors were available for
 389 analysis: gender of infant (1 for male; 0 for female), ethnicity (1 for African American; 0
 390 otherwise), previous pregnancies (1 for yes; 0 for no), cocaine use by mother (1 for yes; 0
 391 for no), age of mother (years), gestational age of infant (weeks), and day of visit during the
 392 study period (a proxy for time since entering the study). The number of visits for each infant
 393 ranged from $m = 2$ to $m = 30$, with a median of $m = 10$ visits. The goal of this analysis was

394 to identify important predictors of infant weight, while accounting for heterogeneity between
 395 infants at baseline and over time.

396 It is natural to include the first four, time-independent covariates (gender, ethnicity, previous
 397 pregnancies, cocaine use) in the model *a-priori* as fixed effects ($p_f = 4$), and to include the
 398 three other time-varying covariates (age of mother, gestational age, day of visit) as composite
 399 effects ($p_c = 3$). An intercept was also included in the model as a composite effect. Prior to
 400 analysis, the three continuous covariates were standardized to have mean zero and variance
 401 one. Adaptive weights were constructed by fitting the saturated model and setting $\nu = 2$.
 402 Using BIC_λ to select the tuning parameter, the final model based on the CREPE estimator
 403 had the following structure

$$\begin{aligned} \hat{\mu}_{ij} &= 6.962 - 0.190 \times \text{gender}_i + 0.245 \times \text{cocaine use}_i + 0.539 \times \text{gestational age}_{ij} \\ &\quad + 2.642 \times \text{visit}_{ij} + b_{0i} + b_i \times \text{visit}_{ij}; \\ \hat{\mathbf{D}} &= \begin{pmatrix} 0.548 & 0.277 \\ 0.277 & 0.214 \end{pmatrix}; \quad \hat{\sigma}^2 = 0.517. \end{aligned}$$

404 Of the four baseline covariates, CREPE identified gender and cocaine dependency as sig-
 405 nificant predictors of infant weight. In particular, prenatal cocaine exposure (PCE) was
 406 associated with higher infant weight, a surprising result given studies previously have found
 407 significant evidence relating PCE and low birth weight (e.g. see the meta-analysis by Gouin
 408 et al. (2011)). Of the time-varying covariates, CREPE identified gestational age as an impor-
 409 tant fixed effect only, and day of visit as an important composite effect, with larger values of
 410 both leading to higher overall infant weights. There was also significant variability between
 411 infant weights at baseline as reflected in the inclusion of a random intercept, in addition to
 412 the variability regarding how weights changed as a function of the day of visit.

413 Comparing the model chosen by CREPE to the one selected using the SCAD and IC_Q method

414 of Ibrahim et al. (2011) (see their Table 2), we find that the latter identified gestational age as
415 (also) having an important random effect, and the age of the mother as having a significant
416 random but not fixed effect, an example of non-hierarchical shrinkage. However, Ibrahim
417 et al. (2011) did not include a random intercept as a candidate covariate, while in our analysis
418 there was substantial variation between infants in their weights at baseline. It is of interest to
419 point out that had we started with the saturated model and applied backwards elimination
420 based on likelihood ratio tests (using `anova` with `lmer` in the R package), then this approach
421 would have produced the same set of informative fixed and random effects as the model
422 selected using CREPE.

423 **7 Discussion**

424 One avenue of research is to extend CREPE to ultra high-dimensional GLMMs, where the
425 number of fixed and/or random effect potentially grows at a faster rate than the number
426 of clusters and cluster size. Such an extension though is of more theoretical interest than
427 of practical relevance. This extension is by no means straightforward: the adaptive weights
428 require modification since the saturated GLMM can no longer be fitted using maximum
429 likelihood estimation (e.g., weights might be constructed based on marginal models, Huang
430 et al. (2008)), and the asymptotic theory demands growing n and m , differing assumptions
431 on the degree of sparsity, and careful consideration of the differing impacts fixed and random
432 effects have on the mixed model.

433 **Supplementary Materials**

434 The proof of Theorem 2, additional simulations results for Gaussian and Bernoulli GLMMs,
435 full results for Poisson GLMMs, and R for implementing the CREPE penalty may be found
436 in the Supplementary Material.

437 **Acknowledgements**

438 This research was supported by the Australian Research Council discovery project grant
439 DP140101259. We are grateful to Andreas Groll for useful discussions.

440 **References**

- 441 Bates, D., Maechler, M., Bolker, B., and Walker, S. (2014). *lme4: Linear mixed-effects*
442 *models using Eigen and S4*. R package version 1.0-6.
- 443 Bondell, H. D., Krishna, A., and Ghosh, S. K. (2010). Joint variable selection for fixed and
444 random effects in linear mixed-effects models. *Biometrics* **66**, 1069–1077.
- 445 Chen, J. and Chen, Z. (2012). Extended BIC for small-n-large-P sparse GLM. *Statistica*
446 *Sinica* **22**, 555–574.
- 447 Cheng, J., Edwards, L. J., Maldonado-Molina, M. M., Komro, K. A., and Muller, K. E.
448 (2010). Real longitudinal data analysis for real people: building a good enough mixed
449 model. *Statistics in Medicine* **29**, 504–520.
- 450 Demidenko, E. (2004). *Mixed Models: Theory and Applications*. Wiley.

- 451 Fan, J. and Li, R. (2001). Variable selection via nonconcave penalized likelihood and its
452 oracle properties. *Journal of the American Statistical Association* **96**, 1348–1360.
- 453 Fan, J. and Peng, H. (2004). Nonconcave penalized likelihood with a diverging number of
454 parameters. *The Annals of Statistics* **32**, 928–961.
- 455 Fan, Y. and Li, R. (2012). Variable selection in linear mixed effects models. *The Annals of*
456 *statistics* **40**, 2043–2068.
- 457 Garcia, T. P., Müller, S., Carroll, R. J., and Walzem, R. L. (2014). Identification of important
458 regressor groups, subgroups and individuals via regularization methods: application to gut
459 microbiome data. *Bioinformatics* **30**, 831–837.
- 460 Gouin, K., Murphy, K., and Shah, P. S. (2011). Effects of cocaine use during pregnancy
461 on low birthweight and preterm birth: systematic review and metaanalyses. *American*
462 *Journal of Obstetrics and Gynecology* **204**, 340.e1 – 340.e12.
- 463 Groll, A. and Tutz, G. (2014). Variable selection for generalized linear mixed models by
464 L_1 -penalized estimation. *Statistics and Computing* **24**, 137–154.
- 465 Hao, X., Yujun, S., Xinjie, W., Jin, W., and Yao, F. (2015). Linear mixed-effects models to
466 describe individual tree crown width for China-Fir in Fujian province, southeast China.
467 *PloS one*, 10:e0122257.
- 468 Huang, J., Ma, S., and Zhang, C. (2008). Adaptive Lasso for sparse high-dimensional
469 regression models. *Statistica Sinica* **18**, 1603–1618.
- 470 Hui, F. K., Warton, D. I., and Foster, S. D. (2015a). Multi-species distribution modeling
471 using penalized mixture of regressions. *The Annals of Applied Statistics* **9**, 866–882.

- 472 Hui, F. K. C., Warton, D. I., and Foster, S. D. (2015b). Tuning parameter selection for the
473 adaptive lasso using ERIC. *Journal of the American Statistical Association* **110**, 262–269.
- 474 Hunter, D. R. and Li, R. (2005). Variable selection using MM algorithms. *The Annals of*
475 *Statistics* **33**, 1617–1642.
- 476 Ibrahim, J. G., Zhu, H., Garcia, R. I., and Guo, R. (2011). Fixed and random effects selection
477 in mixed effects models. *Biometrics* **67**, 495–503.
- 478 Jiang, J., Rao, J. S., Gu, Z., and Nguyen, T. (2008). Fence methods for mixed model
479 selection. *The Annals of Statistics* **36**, 1669–1692.
- 480 Lin, B., Pang, Z., and Jiang, J. (2013). Fixed and random effects selection by REML and
481 pathwise coordinate optimization. *Journal of Computational and Graphical Statistics* **22**,
482 341–355.
- 483 Müller, S., Scealy, J. L., and Welsh, A. H. (2013). Model selection in linear mixed models.
484 *Statistical Science* **28**, 135–167.
- 485 Pan, J. and Huang, C. (2014). Random effects selection in generalized linear mixed models
486 via shrinkage penalty function. *Statistics and Computing* **24**, 725–738.
- 487 Peng, H. and Lu, Y. (2012). Model selection in linear mixed effect models. *Journal of*
488 *Multivariate Analysis* **109**, 109–129.
- 489 Tibshirani, R. (1996). Regression shrinkage and selection via the lasso. *Journal of the Royal*
490 *Statistical Society* **58**, 267–288.
- 491 Tibshirani, R., Saunders, M., Rosset, S., Zhu, J., and Knight, K. (2005). Sparsity and
492 smoothness via the fused lasso. *Journal of the Royal Statistical Society* **67**, 91–108.

- 493 Wasserman, D. and Leventhal, J. (1993). Maltreatment of children born to cocaine-
494 dependent mothers. *American Journal of Diseases of Children* **147**, 1324–1328.
- 495 Wei, G. C. and Tanner, M. A. (1990). A Monte Carlo implementation of the EM algorithm
496 and the poor man’s data augmentation algorithms. *Journal of the American Statistical*
497 *Association* **85**, 699–704.
- 498 Yuan, M. and Lin, Y. (2006). Model selection and estimation in regression with grouped
499 variables. *Journal of the Royal Statistical Society* **68**, 49–67.
- 500 Zhao, P., Rocha, G., and Yu, B. (2009). The composite absolute penalties family for grouped
501 and hierarchical variable selection. *The Annals of Statistics* **37**, 3468–3497.
- 502 Zou, H. (2006). The adaptive lasso and its oracle properties. *Journal of the American*
503 *Statistical Association* **101**, 1418–1429.