

# 1 Joint Selection in Mixed Models using Regularized 2 PQL

3 Francis K.C. Hui<sup>\*1</sup>, Samuel Müller<sup>2</sup>, and A.H. Welsh<sup>1</sup>

4 <sup>1</sup>Mathematical Sciences Institute, The Australian National University,  
5 Canberra, Australia

6 <sup>2</sup>School of Mathematics and Statistics, University of Sydney, Sydney,  
7 Australia

## 8 **Abstract**

9 The application of generalized linear mixed models present some major challenges  
10 for both estimation, due to the intractable marginal likelihood, and model selection, as  
11 we usually want to jointly select over both fixed and random effects. We propose to  
12 overcome these challenges by combining penalized quasi-likelihood (PQL) estimation  
13 with sparsity inducing penalties on the fixed and random coefficients. The resulting  
14 approach, referred to as regularized PQL, is a computationally efficient method for  
15 performing joint selection in mixed models. A key aspect of regularized PQL involves  
16 the use of a group based penalty for the random effects: sparsity is induced such

---

\*Corresponding author: fhui28@gmail.com; Mathematical Sciences Institute, The Australian National University, 0200, Canberra, ACT, Australia

17 that all the coefficients for a random effect are shrunk to zero simultaneously, which  
18 in turns leads to the random effect being removed from the model. Despite being a  
19 quasi-likelihood approach, we show that regularized PQL is selection consistent, i.e.  
20 it asymptotically selects the true set of fixed and random effects, in the setting where  
21 the cluster size grows with the number of clusters. Furthermore, we propose an infor-  
22 mation criterion for choosing the single tuning parameter and show that it facilitates  
23 selection consistency. Simulations demonstrate regularized PQL outperforms [several](#)  
24 [currently employed methods](#) for joint selection even if the cluster size is small com-  
25 pared to the number of clusters, while also offering dramatic reductions in computation  
26 time.

27 **Keywords:** fixed effects, generalized linear mixed models, lasso, penalized likeli-  
28 hood, quasi-likelihood, variable selection

## 29 **1 Introduction**

30 Generalized linear mixed models (GLMMs) are a powerful class of models for analyzing  
31 correlated, non-normal data. Like all regression problems however, model selection is a  
32 difficult but critical part of inference. The problem is especially difficult for mixed models  
33 for two reasons: 1) fitting these models is computationally challenging, and 2) we often  
34 want to jointly select over both the fixed and random effects. Regarding the first problem,  
35 the marginal likelihood for a GLMM has no analytic form except with normal responses  
36 and the identity link, and so numerous estimation methods exist to overcome this diffi-  
37 culty. These range from approximation methods such as penalized quasi-likelihood (PQL,  
38 Breslow and Clayton, 1993), Laplace’s method (Tierney and Kadane, 1986), and numer-  
39 ical quadrature (Rabe-Hesketh et al., 2002), to exact methods such as the Expectation-  
40 Maximization algorithm (EM algorithm, McCulloch, 1997). Of these approaches, PQL

41 is the simplest to implement, as it effectively treats the random effects as “fixed” and es-  
42 timates them in a similar manner to other fixed effects as in a generalized linear model  
43 (GLM). Furthermore, when the cluster size grows with the number of clusters, PQL esti-  
44 mates have been shown to be estimation consistent (Vonesh et al., 2002).

45 For jointly selecting fixed and random effects in GLMMs, proposed methods range  
46 from modifications of information criteria (e.g., Vaida and Blanchard, 2005) to more recent  
47 advances such as the fence (Jiang et al., 2008); see Müller et al. (2013) for an overview  
48 of model selection for LMMs specifically. These methods however are computationally  
49 burdensome to implement, especially since the number of candidate models in the GLMM  
50 context is considerably larger than the GLM context when performing joint selection. One  
51 appealing approach is to use penalized likelihood methods, although their application to  
52 mixed models has only recently been explored. For LMMs, Bondell et al. (2010) pro-  
53 posed adaptive lasso penalties for selecting the fixed and random effects, while Peng and  
54 Lu (2012) and Lin et al. (2013) proposed two-stage methods that separate out the fixed  
55 and random effect selection. For GLMMs, Ibrahim et al. (2011) proposed a modified ver-  
56 sion of the penalty in Bondell et al. (2010), and employed a Monte Carlo EM algorithm  
57 for estimation. This approach however is computationally intensive, with Ibrahim et al.  
58 (2011) limiting their simulations to LMMs only. Focusing solely on computational as-  
59 pects, Schelldorfer et al. (2014) and Groll and Tutz (2014) proposed algorithms for fixed  
60 effects selection only using the lasso penalty in high-dimensional GLMMs, while Pan and  
61 Huang (2014) investigated random effects selection only. The large sample properties of  
62 these algorithms however remain to be determined.

63 In this article, we propose a new approach to joint selection in GLMMs using regular-  
64 ized PQL estimation, and a method to choose the associated tuning parameter. Rather than  
65 working with the marginal likelihood, we propose combining the PQL with adaptive lasso

66 and adaptive group lasso penalties to select the fixed and random effect coefficients re-  
67 spectively. The group lasso is used to exploit the grouped structure inherent in the random  
68 effects: for any random intercept or slope, the coefficients across all clusters are shrunk  
69 to zero at the same time, which leads to the corresponding row and column of the random  
70 effect covariance matrix being shrunk to zero. Such a group penalty approach to random  
71 effects selection has been used previously in linear mixed models by Fan and Li (2012), but  
72 this article is the first to apply it to GLMMs by regularizing the PQL. Another difference  
73 between this article and Fan and Li (2012) is that the latter separate the fixed and random  
74 effects selection into two stages, with different likelihoods and tuning parameters at each  
75 stage, whereas we perform fixed and random effects selection simultaneously using a sin-  
76 gle tuning parameter. Compared to the Monte Carlo EM method of Ibrahim et al. (2011),  
77 joint selection using regularized PQL is extremely fast: it can be viewed as a specific type  
78 of penalized GLM, and the full regularization path can be constructed without the need for  
79 integration.

80 In the setting where the cluster size grows at a slower rate than the number of clusters,  
81 we show that the regularized PQL estimates are estimation and selection consistent. This is  
82 an important advance on Vonesh et al. (2002). For the critical choice of the tuning param-  
83 eter, we propose a new information criterion which we show leads to selection consistency.  
84 This information criterion combines a BIC-type penalty for the fixed effects with an AIC-  
85 type penalty for the random effects. Over the past decade, numerous BIC-type criteria  
86 have been proposed for choosing the tuning parameter in penalized GLMs, particularly  
87 in the high-dimensional setting, with results establishing their selection consistency (e.g.,  
88 Zhang et al., 2010; Hui et al., 2015). Analogous results however do not exist for mixed  
89 models, with the exception of Ibrahim et al. (2011) whose proposed approach involves at  
90 least two tuning parameters. A key contribution of this article is showing that in the case

91 of regularized PQL, differential penalization of the fixed and random effects is needed to  
92 achieve selection consistency.

93 For many applications where the cluster size is small, we propose a hybrid estimator to  
94 improve finite sample performance, i.e. regularized PQL is used for model selection only,  
95 and the final submodel is estimated using maximum likelihood. Simulations demonstrate  
96 that regularized PQL, in conjunction with the proposed information criterion, outperforms  
97 [several currently available methods](#) for joint selection in GLMMs, while offering dramatic  
98 reductions in computation time. We illustrate the application of regularized PQL estimation  
99 on a longitudinal dataset for determining the predictors of forest health over time.

100 To summarize, the main contributions of this article are as follows: 1) we propose a  
101 computationally efficient method of performing joint selection in GLMMs, which com-  
102 bines the PQL with adaptive (group) lasso penalties to regularize the fixed and random  
103 effect coefficients; 2) we develop an information criterion for choosing the tuning param-  
104 eter in regularized PQL estimation, that involves differing model complexity terms on the  
105 fixed and random effects; 3) we demonstrate estimation and selection consistency prop-  
106 erties for regularized PQL estimation, and show that the proposed information criterion  
107 asymptotically chooses a tuning parameter that leads to selection consistency; 4) we per-  
108 form simulations to demonstrate the computational speed and strong performance of regu-  
109 larized PQL, relative to other penalized likelihood methods, even when the cluster size is  
110 relatively small compared to the number of clusters.

## 111 **2 Generalized Linear Mixed Models**

112 We focus on the independent cluster model with random intercepts and slopes. Let  $y_{ij}$   
113 denote the  $j^{\text{th}}$  measurement for the  $i^{\text{th}}$  cluster, where  $i = 1, \dots, n$  and  $j = 1, \dots, m_i$ .

114 Note that we allow for unequal cluster sizes. Let  $\mathbf{x}_{ij}$  be a vector of  $p_f$  covariates corre-  
 115 sponding to fixed effects, and  $\mathbf{z}_{ij}$  be a vector of  $p_r$  covariates corresponding to random  
 116 effects. Both  $\mathbf{x}_{ij}$  and  $\mathbf{z}_{ij}$  may contain an intercept term as their first element. We as-  
 117 sume that  $p_f$  and  $p_r$  are fixed, with  $p_r < \min_i(m_i)$  where  $\min_i(\cdot)$  denotes the minimum  
 118 over  $i = 1, \dots, n$ . Conditional on the random effects  $\mathbf{b}_i$ , the responses  $y_{ij}$  are assumed to  
 119 come from a distribution in the exponential family, with density function  $f(y_{ij}|\boldsymbol{\beta}, \mathbf{b}_i, \phi) =$   
 120  $\exp[\{y_{ij}\vartheta_{ij} - a(\vartheta_{ij})\}/\phi + c(y_{ij}, \phi)]$  for known functions  $a(\cdot)$  and  $c(\cdot)$  and dispersion pa-  
 121 rameter  $\phi$ . The mean,  $\mu_{ij}$ , is modeled as  $g(\mu_{ij}) = \eta_{ij} = \mathbf{x}_{ij}^T \boldsymbol{\beta} + \mathbf{z}_{ij}^T \mathbf{b}_i$ , for a known link  
 122 function  $g(\cdot)$ . For simplicity, we assume the canonical link is used, so  $g(\mu_{ij}) = \vartheta_{ij} = \eta_{ij}$   
 123 and  $\mu_{ij} = a'(\eta_{ij})$ . The random effects are normally distributed,  $\mathbf{b}_i \sim \mathcal{N}_{p_r}(\mathbf{0}, \mathbf{D})$ , where  $\mathbf{D}$   
 124 is the random effect covariance matrix.

125 For the  $i^{\text{th}}$  cluster, we have an  $m_i$ -vector  $\mathbf{y}_i = (y_{i1}, \dots, y_{im_i})$ , a  $m_i \times p_f$  matrix  $\mathbf{X}_i =$   
 126  $(\mathbf{x}_{i1} \dots \mathbf{x}_{im_i})^T$  of fixed effect covariates, and a  $m_i \times p_r$  matrix  $\mathbf{Z}_i = (\mathbf{z}_{i1} \dots \mathbf{z}_{im_i})^T$  of  
 127 random effect covariates. In turn, we can write  $g(\boldsymbol{\mu}_i) = \boldsymbol{\eta}_i = \mathbf{X}_i^T \boldsymbol{\beta} + \mathbf{Z}_i^T \mathbf{b}_i$ , where  $g(\cdot)$  is  
 128 applied component-wise,  $\boldsymbol{\beta} = (\beta_1, \dots, \beta_{p_f})$ ,  $\mathbf{b}_i = (b_{i1}, \dots, b_{ip_r})$ ,  $\boldsymbol{\mu}_i = (\mu_{i1}, \dots, \mu_{im_i})$  and  
 129 similarly for  $\boldsymbol{\eta}_i$ . Finally, let  $\mathbf{b} = (\mathbf{b}_1^T, \dots, \mathbf{b}_n^T)^T$  be the  $np_r$ -vector of all random effects, and  
 130  $\boldsymbol{\Psi} = \{\boldsymbol{\beta}^T, \text{vech}(\mathbf{D})^T\}^T$  where  $\text{vech}(\cdot)$  denotes the half-vectorization operator. Note that  
 131 each  $\mathbf{b}_i$  is of fixed dimension  $p_r$ , while  $\dim(\mathbf{b})$  grows with linearly with  $n$ .

132 For the GLMM above, the marginal log-likelihood is

$$\ell(\boldsymbol{\Psi}) = -\frac{n}{2} \log \det(\mathbf{D}) + \sum_{i=1}^n \log \left( \int \exp \left( \sum_{j=1}^{m_i} \log f(y_{ij}|\boldsymbol{\beta}, \mathbf{b}_i) - \frac{1}{2} \mathbf{b}_i^T \mathbf{D}^{-1} \mathbf{b}_i \right) d\mathbf{b}_i \right),$$

133 where  $\det(\mathbf{D})$  is the determinant of  $\mathbf{D}$ . Aside from linear mixed models, the integral in  
 134 the marginal log-likelihood does not have an analytic form, and this complicates maximum  
 135 likelihood estimation. A popular, alternative estimation method is PQL estimation, which

136 involves maximizing the quasi-likelihood function

$$\ell_{\text{PQL}}(\Psi, \mathbf{b}) = \sum_{i=1}^n \sum_{j=1}^{m_i} \log f(y_{ij} | \beta, \mathbf{b}_i) - \frac{1}{2} \sum_{i=1}^n \mathbf{b}_i^T \mathbf{D}^- \mathbf{b}_i, \quad (1)$$

137 where  $\mathbf{D}^-$  denotes the Moore-Penrose generalized inverse of  $\mathbf{D}$ . The use of a generalized  
 138 inverse here, as opposed to the standard matrix inverse, allows us to deal with cases where  
 139 the covariance matrix is singular (see Breslow and Clayton, 1993). This is necessary when  
 140 we establish asymptotic properties in Section 4, where the true random effects are assumed  
 141 to be sparse.

142 There is a close link between PQL estimation and Laplace's method for GLMMs.  
 143 Specifically, for a fixed  $\Psi$ , let  $\tilde{\mathbf{b}} = (\tilde{\mathbf{b}}_1^T, \dots, \tilde{\mathbf{b}}_n^T)^T$  denote the maximizer of (1). Then  
 144 the Laplace approximated log-likelihood is defined as

$$\ell_{\text{LA}}(\Psi) = \sum_{i=1}^n \sum_{j=1}^{m_i} \log f(y_{ij} | \beta, \tilde{\mathbf{b}}_i) - \frac{1}{2} \sum_{i=1}^n \tilde{\mathbf{b}}_i^T \mathbf{D}^- \tilde{\mathbf{b}}_i - \frac{1}{2} \sum_{i=1}^n \log \det(\mathbf{Z}_i^T \tilde{\mathbf{W}}_i \mathbf{Z}_i \mathbf{D} + \mathbf{I}_{p_r}),$$

145 where  $\tilde{\mathbf{W}}_i$  is a  $m_i \times m_i$  diagonal weight matrix with elements  $\tilde{\mathbf{W}}_{i,jj} = a''(\tilde{\eta}_{ij})/\phi$ ,  $\tilde{\eta}_{ij} =$   
 146  $\mathbf{x}_{ij}^T \beta + \mathbf{z}_{ij}^T \tilde{\mathbf{b}}_i$ , and  $\mathbf{I}_{p_r}$  is an identity matrix of dimension  $p_r$ . The key difference between  
 147 PQL and the Laplace approximation lies in the last term, which is a non-linear function of  
 148  $\beta$  and  $\mathbf{b}$ . By assuming the weights in  $\mathbf{W}_i$  vary slowly with the mean, Breslow and Clayton  
 149 (1993) proposed ignoring this last term, from which the PQL follows. Note that when [the](#)  
 150 [minimum cluster size](#)  $\min_i(m_i)$ , and hence all  $m_i$ , are large, the estimates from PQL and  
 151 Laplace's method should be close to each other, since the last term in  $\ell_{\text{LA}}(\Psi)$  is of a smaller  
 152 order than the first term (see also Demidenko, 2013, Section 7.3). For normal responses,  
 153  $a''(\eta_{ij}) = 1$ , so the estimates of  $\beta$  based on  $\ell_{\text{LA}}(\Psi, \mathbf{b})$  and  $\ell_{\text{PQL}}(\Psi, \mathbf{b})$  coincide, noting that  
 154 the Laplace approximation is exact for normal linear mixed models.

155 Compared to maximizing the marginal and Laplace approximated log-likelihoods, PQL  
 156 estimation is straightforward: equation (1) resembles the log-likelihood for a GLM com-  
 157 bined with a generalized ridge penalty, where  $\mathbf{b}$  is also treated as a fixed effect vector,  
 158 and so modifications of standard optimization routines such as iteratively reweighted least  
 159 squares can be used for maximization. This in turn motivates us to consider using the PQL  
 160 as a loss function for penalized joint selection in GLMMs.

### 161 3 Regularized PQL Estimation

162 We propose regularized PQL estimation to perform selection over both the fixed and ran-  
 163 dom effects in GLMMs.

164 **Definition.** For a given  $\mathbf{D}$ , the regularized PQL estimates of the fixed and random effect  
 165 coefficients are given by

$$(\hat{\boldsymbol{\beta}}_\lambda, \hat{\mathbf{b}}_\lambda) = \arg \max_{\boldsymbol{\beta}, \mathbf{b}} \ell_p(\boldsymbol{\Psi}, \mathbf{b}) = \arg \max_{\boldsymbol{\beta}, \mathbf{b}} \ell_{PQL}(\boldsymbol{\Psi}, \mathbf{b}) - \lambda \sum_{k=1}^{p_f} v_k |\beta_k| - \lambda \sum_{l=1}^{p_r} w_l \|\mathbf{b}_{\bullet l}\|,$$

166 where  $v_k$  and  $w_k$  are adaptive weights based on preliminary estimates of  $\beta_k$  and  $\mathbf{D}$  respec-  
 167 tively,  $\mathbf{b}_{\bullet l} = (b_{1l}, \dots, b_{nl})$  denotes all the coefficients corresponding to the  $l^{\text{th}}$  random effect,  
 168 and  $\|\cdot\|$  denotes its  $L_2$  norm.

169 We use an adaptive lasso penalty with weights  $v_k$  for the fixed effects, and an adaptive  
 170 group lasso penalty with weights  $w_l$  for the random effects, linked by one tuning parameter  
 171  $\lambda > 0$ . Specifically, let  $\tilde{\boldsymbol{\beta}}$  and  $\tilde{\mathbf{D}}$  denote the unpenalized, maximum likelihood estimates of  
 172 the fixed effect coefficients and random effect covariance matrix respectively from fitting  
 173 the full GLMM. This fitting could be performed, for example, by applying the EM algo-  
 174 rithm, or via recent advances in maximum likelihood estimation for GLMMs such as data



175 cloning (Lele et al., 2010). Then we choose  $v_k = |\tilde{\beta}_k|^{-\kappa}$  and  $w_l = \tilde{D}_{ll}^{-\kappa}$ , where  $\tilde{D}_{ll}$  is the  
 176  $l^{\text{th}}$  diagonal element of  $\tilde{D}$  and  $\kappa > 0$  is a common power parameter. Note that while the  
 177 penalty involves  $\mathbf{b}$ , the adaptive weights for the random effects require only an initial esti-  
 178 mate of  $\mathbf{D}$ . Also, in the case where the fixed intercept term is included but not penalized,  
 179 the adaptive lasso penalty is summed from  $k = 2$  to  $p_f$ .

180 The adaptive weights mean that a single tuning parameter, as opposed to using different  
 181  $\lambda$ 's for the fixed and random effects, is able to achieve consistency of the regularized PQL  
 182 estimates. In Section 3.2, we discuss how to select the tuning parameter. Of course, having  
 183 to select over multiple  $\lambda$ 's also presents a considerable computational challenge (see for  
 184 instance, Garcia et al., 2014). Note that due to the concavity of both  $\ell_{\text{PQL}}(\Psi, \mathbf{b})$  and the  
 185 lasso penalties, if there exists a maximizer to  $\ell_p(\Psi, \mathbf{b})$  then it is also the unique, regularized  
 186 PQL estimate (see also Lemma 2.1, Jiang et al., 2001).

187 Regularized PQL performs joint selection of the fixed and random effects in mixed  
 188 models. The adaptive group lasso penalizes random slopes across clusters, thereby utiliz-  
 189 ing the grouped structure inherent in the random effects. For a sufficiently large value of  $\lambda$ ,  
 190 maximizing the regularized PQL shrinks  $\|\mathbf{b}_{\bullet l}\| = 0$ , that is, all the coefficients correspond-  
 191 ing to the  $l^{\text{th}}$  random slope (or the random intercept) are shrunk to zero. This implies that  
 192 the  $l^{\text{th}}$  row and column of  $\mathbf{D}$  are also set to zero (see Section 3.1). This method of penaliz-  
 193 ing the coefficients  $\mathbf{b}$  explicitly differs from the random effects penalties that shrink one or  
 194 more elements of  $\mathbf{D}$  or a decomposition of  $\mathbf{D}$  to zero (Bondell et al., 2010; Ibrahim et al.,  
 195 2011). In fact, the potential to penalize  $\mathbf{b}$  arises precisely because the PQL is a function of  
 196 the  $\mathbf{b}$ 's.

197 Since  $\ell_p(\Psi, \mathbf{b})$  does not require integrating over the random effects, the solution path  
 198 for the regularized PQL estimates is easily constructed. Conditional on  $\mathbf{D}$  and  $\mathbf{b}$ , estimates  
 199 of the fixed effects  $\beta$  are obtained by fitting a GLM with the adaptive lasso penalty across

200 all clusters, with  $\mathbf{z}_{ij}^T \mathbf{b}_i$  as an offset. Then conditional on  $\mathbf{D}$  and  $\beta$ , estimates of the random  
 201 effects  $\mathbf{b}$  are obtained by fitting a GLM with an adaptive elastic net penalty, with  $\mathbf{x}_{ij}^T \beta_i$  as  
 202 an offset. In the simulations and application, we used a local quadratic approximation (Fan  
 203 and Li, 2001) to calculate the regularized PQL estimates, and this was already consider-  
 204 ably faster than methods involving the marginal likelihood. Utilizing more sophisticated  
 205 methods for estimation (e.g., coordinate descent, Friedman et al., 2010) will further reduce  
 206 computation time.

### 207 3.1 Estimation of the Covariance Matrix

208 For a given  $\mathbf{D}$ , regularized PQL provides estimates of the fixed and random effect coeffi-  
 209 cients  $(\hat{\beta}_\lambda^T, \hat{\mathbf{b}}_\lambda^T)^T$ . With these estimates, we can update the random effect covariance matrix  
 210 in a number of ways (e.g., Breslow and Clayton, 1993; Vonesh et al., 2002). We propose  
 211 substituting  $(\hat{\beta}_\lambda^T, \hat{\mathbf{b}}_\lambda^T)^T$  back into  $\ell_{\text{LA}}(\Psi)$ , and then maximizing to obtain an estimate of  $\mathbf{D}$ .  
 212 Straightforward algebra (see Appendix A) shows that an estimate of the covariance matrix  
 213 can be obtained via the following iterative equation: At the  $t^{\text{th}}$  iteration,

$$\hat{\mathbf{D}}_\lambda^{(t)} = \frac{1}{n} \sum_{i=1}^n \left\{ \left( \mathbf{z}_i^T \hat{\mathbf{W}}_{\lambda i} \mathbf{z}_i + (\hat{\mathbf{D}}_\lambda^{(t-1)})^{-1} \right)^{-1} + \hat{\mathbf{b}}_{\lambda i} \hat{\mathbf{b}}_{\lambda i}^T \right\}, \quad (2)$$

214 where  $\hat{\mathbf{b}}_\lambda^T = (\hat{\mathbf{b}}_{\lambda 1}^T, \dots, \hat{\mathbf{b}}_{\lambda n}^T)^T$  and  $\hat{\mathbf{W}}_{\lambda i}$  is the weight matrix for subject  $i$  evaluated at  
 215  $(\hat{\beta}_\lambda^T, \hat{\mathbf{b}}_{\lambda i}^T)^T$ . Note that when  $\|\mathbf{b}_{\bullet i}\|$  is shrunk to zero, it makes sense to set the  $i^{\text{th}}$  row and  
 216 column of  $\mathbf{D}$  to zero, reflecting the removal of this covariate from the random effects com-  
 217 ponent of the model. In such a case, the iterative formula above is applied only to the  
 218 submatrix of  $\mathbf{D}$  with non-zero rows and columns. Finally, we point out that this update  
 219 of the covariance matrix is only used in the context of regularized PQL estimation; as we  
 220 discuss in Section 3.3, we propose using a hybrid estimator to calculate the final parameter

221 estimates.

## 222 3.2 Tuning Parameter Selection

223 As with all penalized likelihood methods, both the finite sample and asymptotic perfor-  
224 mance of regularized PQL depend critically on being able to choose an appropriate value  
225 of the tuning parameter. For the GLM framework, there has been considerable research  
226 into choosing  $\lambda$  using, most commonly, cross validation or information criteria (e.g., Zhang  
227 et al., 2010), and we focus on the latter method. Specifically, we consider tuning parameters  
228 within the range  $[\lambda_{\min}, \lambda_{\max}]$ , where  $\lambda_{\min}$  leads to the full model containing all the candi-  
229 date fixed and random effects, and  $\lambda_{\max}$  is the smallest  $\lambda$  that leads to the null model. A  
230 solution path is constructed by considering a sequence of  $\lambda$ 's over this range, and selecting  
231 the value of  $\lambda$  (hence the best submodel) by minimizing the [information criterion](#)

$$\text{IC}(\lambda) = -\frac{2}{N}\ell_{\text{PQL}}(\hat{\Psi}_\lambda, \hat{\mathbf{b}}_\lambda) + \frac{\log(n)}{N} \dim(\hat{\beta}_\lambda) + \frac{2}{N} \dim(\hat{\mathbf{b}}_\lambda), \quad (3)$$

232 where  $\dim(\hat{\beta}_\lambda)$  and  $\dim(\hat{\mathbf{b}}_\lambda)$  are the number of non-zero estimated fixed and random effect  
233 coefficients respectively and, importantly,  $\dim(\hat{\mathbf{b}}_\lambda) = n \dim(\hat{\mathbf{b}}_{\lambda 1})$ . [Note that division by](#)  
234 [total sample size  \$N\$  is often used when studying information criteria for tuning parameter](#)  
235 [selection \(e.g., Zhang et al., 2010\).](#)

236 A key feature of  $\text{IC}(\lambda)$ , which sets it apart from standard information criteria used for  
237 tuning parameter selection in other penalties for mixed models (e.g., Ibrahim et al., 2011;  
238 Lin et al., 2013), is its use of different model complexity penalties. Specifically, a BIC-type  
239 penalty of ' $\log(n)$ ' is used for the fixed effects, and an AIC-type penalty of '2' is used for  
240 the random effects. The latter arises because the model complexity is already taken into  
241 account by  $\dim(\hat{\mathbf{b}}_\lambda)$ , which grows linearly with  $n$  regardless of the number of random ef-

242 facts in the model. Put another way, overfitting of the  $\hat{\mathbf{b}}_\lambda$ 's is inherently prevented by the  
 243 information criterion, since the removal of one random effect from the model amounts to  
 244 the removal of  $n$  coefficients in regularized PQL by the group sparsity of  $\|\hat{\mathbf{b}}_{\lambda \bullet i}\|$ . By con-  
 245 trast,  $\dim(\hat{\boldsymbol{\beta}}_\lambda)$  is always of order  $O_p(1)$ , and so the BIC-type penalty of  $\log(n)$  is necessary  
 246 to properly account for model complexity in the fixed effects and prevent overfitting (see  
 247 Shao, 1997, for related work on the use of differing model complexities in the linear re-  
 248 gression context). In Section 4.1, we show that using  $\text{IC}(\lambda)$  to choose the tuning parameter  
 249 selection leads to selection consistency in regularized PQL.

### 250 3.3 Hybrid Estimation Approach

251 In real finite sample settings, regularized PQL can produce biased estimates of the fixed  
 252 effects and the random effect covariance matrix. The bias is related to the well known finite  
 253 sample bias for unpenalized PQL estimation [when the cluster sizes are not](#) large compared  
 254 to the number of clusters (Lin and Breslow, 1996). Moreover, as shown in Theorem 1, we  
 255 establish [consistency of the regularized PQL estimates where the convergence rate depends](#)  
 256 [on the rate of growth of the cluster sizes  \$m\_i\$](#) . Thus compared to the unpenalized maximum  
 257 likelihood estimates, which are  $n^{1/2}$ -consistent, [the regularized PQL estimates are not as](#)  
 258 [efficient if all the  \$m\_i\$ 's are smaller than  \$n\$ , which is typically the case with longitudinal](#)  
 259 [studies](#).

260 To improve finite sample performance, we propose a hybrid estimation approach in  
 261 which we use regularized PQL *only* for joint selection of the fixed and random effects, and  
 262 use maximum likelihood estimation of the selected submodel to obtain the final estimates  
 263  $\boldsymbol{\beta}$  and  $\text{vech}(\mathbf{D})$ , [as well as to construct predictions of the random effects based on posterior](#)  
 264 [modes \(for instance\)](#). Hybrid estimation approaches have been used previously (e.g., Hui  
 265 et al., 2015), although the purpose there was to reduce the bias introduced by penalization,

266 while we use the hybrid approach to address both the relative lack of efficiency and finite  
 267 sample bias of the regularized PQL estimates. Of course, since the hybrid approach is  
 268 applied on the submodel chosen by regularized PQL estimation, it also inherits the selection  
 269 consistency property encapsulated in the second part of Theorem 1. In the simulations  
 270 in Section 5, we empirically evaluate the performance of the hybrid estimation approach  
 271 compared to just using the estimates from regularized PQL.

## 272 4 Asymptotic Properties

273 We study the large sample properties of regularized PQL estimation in the setting where the  
 274 cluster sizes grow with the number of clusters. Without loss of generality, suppose that the  
 275 clusters are labeled such that the first cluster grows at the slowest rate, and the last cluster  
 276 grows at the largest rate. That is,  $m_1 = O(m_k)$  for all  $k = 2, \dots, n$ , and  $m_l = O(m_n)$  for  
 277 all  $l = 1, \dots, n_1$ , so the rates of growth of the cluster sizes are bounded below by the order  
 278 of  $m_1$  and above by the order of  $m_n$ . Note this includes the case where all cluster sizes  
 279 are constrained to grow at the same rate. It is also worth pointing out that no restriction  
 280 is made directly on whether the cluster sizes are balanced or not; Instead, the assumptions  
 281 made concern the rate of growth of the cluster sizes. We assume that  $m_n/n \rightarrow 0$ , such that  
 282 all cluster sizes grow at a smaller rate than number of clusters. This setting arises commonly  
 283 in longitudinal studies in epidemiology (for instance), where the number of measurements  
 284 recorded for each cluster increases slowly as more subjects are recruited into the study.

285 To aid our theoretical development, write the random effect covariance matrix as  $\mathbf{D} =$   
 286  $\mathbf{\Gamma}\mathbf{\Gamma}^T$ , where  $\mathbf{\Gamma} = \mathbf{Q}\mathbf{\Lambda}^{1/2}$  with  $\mathbf{Q}$  the orthogonal matrix of normalized eigenvectors and  $\mathbf{\Lambda}$   
 287 the diagonal matrix whose entries are the eigenvalues of  $\mathbf{D}$ . Note if the  $l^{\text{th}}$  row of  $\mathbf{\Gamma}$  is equal  
 288 to zero, then it implies that both the  $l^{\text{th}}$  row and column of  $\mathbf{D}$  are zero. Consequently, for

289 the remainder of this section, we redefine the parameter vector as  $\Psi = \{\beta^T, \text{vec}(\Gamma)^T\}^T \in$   
 290  $\Re^{p_f + p_r^2}$ , replacing  $\text{vech}(\mathbf{D})$  by  $\text{vec}(\Gamma)$ . This parameterization is used only in the theoretical  
 291 development, as it avoids the true parameter point being on the boundary of the parameter  
 292 space (see Condition C4 below) and is not employed in the estimation process.

293 Let  $\Psi_0 = \{\beta_0^T, \text{vec}(\Gamma_0^T)\}^T$  be the true parameter point and, without loss of generality,  
 294 write  $\beta_0 = (\beta_{01}^T, \beta_{02}^T = \mathbf{0}^T)^T$  and  $\text{vec}(\Gamma_0) = (\text{vec}(\Gamma_{01}^T), \text{vec}(\Gamma_{02}^T) = \mathbf{0}^T)^T$ . Let  $p_{0f} =$   
 295  $\dim(\beta_{01})$  denote the number of truly non-zero fixed effects, and  $p_{0r}$  the number of rows  
 296 in  $\Gamma_{01}$ . Also, for  $i = 1, \dots, n$ , let  $\mathbf{b}_{0i}$  denote a realization from the true random effects  
 297 distribution; the first  $p_{0r}$  elements of  $\mathbf{b}_{0i}$  are drawn from a multivariate normal distribution  
 298 with mean zero and covariance matrix  $\mathbf{D}_{01} = \Gamma_{01}\Gamma_{01}^T$ , and  $b_{0il} = 0$  for  $l = p_{0r} + 1, \dots, p_r$ .  
 299 Finally, let  $N = \sum_{i=1}^n m_i$  be the total number of observations. The following regularity  
 300 conditions are required.

301 (C1) The function  $a(\eta)$  is three times continuously differentiable in its domain, with

302 
$$a''(\eta) \geq c_0 > 0 \text{ for some sufficiently small constant } c_0.$$

303 (C2) For every  $i = 1, \dots, n$  and  $j = 1, \dots, m_i$ , there exists a sufficiently large constant

304  $C$  such that  $\|\mathbf{x}_{ij}\|_\infty < C$  and  $\|\mathbf{z}_{ij}\|_\infty < C$  where  $\|\cdot\|_\infty$  is the maximum norm.

305 Furthermore, the matrices  $m_i^{-1} \mathbf{X}_i^T \mathbf{X}_i$  and  $m_i^{-1} \mathbf{Z}_i^T \mathbf{Z}_i$  are positive definite with min-

306 imum and maximum eigenvalues bounded from above and below by  $1/c_1$  and  $c_1$

307 respectively, where  $c_1$  is some positive constant.

308 (C3) Let  $\ell_1(\beta, \mathbf{b}) = \sum_{i=1}^n \sum_{j=1}^{m_i} \log f(y_{ij} | \beta, \mathbf{b}_i)$ , and  $\mathbf{H}(\beta, \mathbf{b}) = -N^{-1} \nabla^2 \ell_1(\beta, \mathbf{b})$ . Then there

309 exists a  $\varepsilon > 0$  such that for  $n$  and all  $m_i$  sufficiently large, the minimum eigenvalue

310 of  $\mathbf{H}(\beta, \mathbf{b})$  is bounded away from zero for all  $\|(\beta^T, \mathbf{b}^T)^T - (\beta_0^T, \mathbf{b}_0^T)^T\|_\infty \leq \varepsilon$ .

311 (C4)  $\Psi_0$  is a interior point in the compact set  $\Omega \in \Re^{p_f + p_r^2}$ .

312 (C5) The tuning parameter  $\lambda$  satisfies (a)  $\lambda m_1^{1/2}/N \rightarrow 0$  and (b)  $\lambda m_1^{1/2} n^{\kappa/2}/N \rightarrow \infty$ ,  
 313 where  $m_n/n \rightarrow 0$ .

314 Conditions (C1) and (C2) ensure the observed information matrices based on  $\ell_{\text{PQL}}(\Psi, \mathbf{b})$   
 315 are positive definite, and imply that the expectations  $E_{\mathbf{b}}\{a''(\eta)\}$  and  $E_{\mathbf{b}}\{a'''(\eta)\}$ , where the  
 316 expectations are with respect to the true random effects distribution, are finite. Condition  
 317 (C3) extends this to a small neighborhood around the true parameters. Along with the  
 318 independence of the responses  $\mathbf{y}_i$  for each cluster, conditions (C1)-(C4) are sufficient to  
 319 ensure the maximum likelihood estimate of  $\Psi$ , i.e. the maximizer  $\ell(\Psi)$ , exists and is  $n^{1/2}$ -  
 320 consistent (Lehmann, 1983). Condition (C5) imposes restriction on the rate at which the  
 321 tuning parameter can grow.

322 We now present a result on the large sample consistency of the regularized PQL esti-  
 323 mates.

324 **Theorem 1.** *Under conditions (C1)-(C5a), as  $n, m_i \rightarrow \infty$  for all  $i$  and  $m_n/n \rightarrow 0$ ,*  
 325 *the regularized PQL estimator satisfies  $\|\hat{\beta}_\lambda - \beta_0\| = O_p(m_1^{-1/2})$  and  $\|\hat{\mathbf{b}}_{\lambda i} - \mathbf{b}_{0i}\| =$*   
 326  *$O_p(m_1^{-1/2})$  for all  $i = 1, \dots, n$ . If condition (C5b) is also satisfied, then  $P(\hat{\beta}_{\lambda 02} =$*   
 327  *$\mathbf{0}) \rightarrow 1$  and  $P(\|\hat{\mathbf{b}}_{\lambda \bullet l}\| = 0) \rightarrow 1$  for all  $l = p_{0r}+1, \dots, p_r$ , where  $\hat{\mathbf{b}}_{\lambda \bullet l} = (\hat{b}_{\lambda 1l}, \dots, \hat{b}_{\lambda nl})$*   
 328 *denotes all the estimated coefficients corresponding to the  $l^{\text{th}}$  random effect.*

329 Note that even though  $p_r$  is fixed, each  $\hat{\mathbf{b}}_{\lambda \bullet l}$  is growing at the same rate as the number  
 330 of clusters,  $n$ , so the proof of Theorem 1 has to be developed in a high-dimensional set-  
 331 ting. Outlines of the proofs of all theorems are given in Appendix B, with detailed proofs  
 332 provided in the Supplementary Material.

333 The  $m_1^{1/2}$ -consistency of the fixed effects agrees with the result of Vonesh et al. (2002),  
 334 who showed  $\hat{\beta}_\lambda - \beta_0 = O_p(\max\{m^{-1/2}, n^{-1/2}\})$  in the case where all cluster sizes were  
 335 equal to  $m$ , and  $n \rightarrow \infty$  and  $m \rightarrow \infty$ . The  $m_1^{1/2}$ -consistency for each  $\hat{\mathbf{b}}_{\lambda i}$  is also reason-

336 able, since regularized PQL treats the  $\mathbf{b}$  as fixed effects and the estimation of  $\mathbf{b}_{\lambda_i}$  depends  
337 only on the  $m_i$  observations within the  $i^{\text{th}}$  cluster. Estimation consistency for all random  
338 effect coefficients is thus governed by the smallest rate of growth of the cluster sizes,  $m_1$ .  
339 The second part of Theorem 1 states that the regularized PQL estimators asymptotically  
340 select only the truly non-zero fixed and random effects in the GLMM. Together with its  
341 computational simplicity, this presents a strong argument for the use of regularized PQL  
342 for joint selection.

## 343 4.1 Consistency of $\text{IC}(\lambda)$

344 In this section, we show that using the tuning parameter chosen by minimizing  $\text{IC}(\lambda)$   
345 asymptotically identifies the true model. For any value of  $\lambda \in [\lambda_{\min}, \lambda_{\max}]$ , let  $\alpha$  denote  
346 the submodel (subset of fixed and random effects) selected by regularized PQL estima-  
347 tion. Clearly  $\alpha$  depends on  $\lambda$ , but for ease of notation we have suppressed this dependence.  
348 Next, let  $(\tilde{\Psi}_\alpha^T, \tilde{\mathbf{b}}_\alpha^T)^T$  denote the unregularized PQL estimator for this submodel, obtained  
349 by maximizing the PQL in (1) along with the iterative update of the covariance matrix in  
350 (2). Finally, let  $\lambda_0$  be a sequence of tuning parameters that satisfy condition (C5) and hence  
351 selects the true model, which we denote here as  $\alpha_0$ .

352 For the development below, we require an additional condition. Let ‘ $\supset$ ’ denote the  
353 proper superset relation.

354 (C6) There exists a constant  $c_2$  such that  $E \{ \ell_o(\Psi_0) - \ell_o(\Psi_\alpha^*) \} \geq c_2 > 0$  for all mod-  
355 els  $\alpha \not\supset \alpha_0$ , where  $\ell_o(\Psi)$  denotes the marginal log-likelihood of a GLMM for a  
356 single observation, and  $\Psi_\alpha^*$  denotes the pseudo-true parameters for model  $\alpha$  which  
357 minimize  $E \{ -\ell_o(\Psi_\alpha^*) \}$ .

358 Conditions like (C6) are imposed in theoretical developments on selection consistency e.g.,



359 see condition (viii) in Müller and Welsh (2009) for robust selection on GLMs, and condi-  
 360 tion (C4) in Zhang et al. (2010) in the setting of penalized GLMs. It amounts to requiring  
 361 that the Kullback-Leibler distance between any underfitted GLMM, with pseudo-true pa-  
 362 rameters  $\Psi_\alpha^*$ , and the true GLMM is positive; see the Supplementary Material and White  
 363 (1982) for further discussion of pseudo-true parameters.

364 We now define a proxy information criterion based on these unregularized PQL esti-  
 365 mates,

$$\text{IC}_{\text{proxy}}(\alpha) = -\frac{2}{N} \sum_{i=1}^n \sum_{j=1}^{m_i} \log f(y_{ij} | \tilde{\beta}_\alpha, \tilde{\mathbf{b}}_{\alpha i}) + \frac{\log(n)}{N} \dim(\tilde{\beta}_\alpha) + \frac{2}{N} \dim(\tilde{\mathbf{b}}_\alpha).$$

366 Note the loss function for this proxy criterion involves only the first part of the PQL. The  
 367 reason for introducing this proxy criterion is to simplify the theoretical development: since  
 368  $\text{IC}_{\text{proxy}}(\alpha)$  does not involve penalized estimates, we can focus on establishing its asymptotic  
 369 behavior when  $\alpha$  represents an underfitted or overfitted model without having to deal with  
 370 the effects of  $\lambda$ . We then have the following result.

371 **Lemma 1.** *Under conditions (C1)-(C4) and (C6), and as  $n, m_i \rightarrow \infty$  for all  $i$  and  $m_n/n \rightarrow$   
 372  $0$ , the proxy information criterion satisfies  $P \{ \min_{\alpha \neq \alpha_0} \text{IC}_{\text{proxy}}(\alpha) > \text{IC}_{\text{proxy}}(\alpha_0) \} \rightarrow 1$ .*

373 Lemma 1 guarantees that asymptotically, all underfitted (at least one truly non-zero  
 374 coefficient is missing from the model) and overfitted (all truly non-zero coefficients and one  
 375 or more zero coefficient are included in the model) models estimated using unregularized  
 376 PQL will have values of  $\text{IC}_{\text{proxy}}(\alpha)$  greater than the value attained at the true model  $\alpha_0$ .  
 377 From these results, we are able to infer the large sample properties of  $\text{IC}(\lambda)$  for choosing  
 378 the tuning parameter.

379 **Theorem 2.** *Let  $\hat{\alpha}$  be the model chosen by minimizing  $\text{IC}(\lambda)$  defined in (3). Then under  
 380 conditions (C1)-(C4) and (C6), and as  $n, m_1 \rightarrow \infty$ , it holds that  $P(\hat{\alpha} = \alpha_0) \rightarrow 1$ .*

381 The above guarantees that the model chosen by minimizing  $\text{IC}(\lambda)$  is asymptotically  
 382 equal to the model chosen by  $\lambda_0$ . Since  $\lambda_0$  satisfies condition (C5) and selects the true  
 383 model, it follows immediately that choosing the tuning parameter based on  $\text{IC}(\lambda)$  leads to  
 384 consistent model selection using regularized PQL.

## 385 5 Simulation Study

386 We performed an empirical study to assess the performance of regularized PQL estima-  
 387 tion and  $\text{IC}(\lambda)$  for three commonly applied forms of GLMMs, namely the linear mixed  
 388 model, Bernoulli, and Poisson GLMMs. For brevity, we only present the first two sets  
 389 of results here; the Poisson GLMM results are presented in the Supplementary Material.  
 390 For simplicity, we restrict our simulations to cases where the cluster sizes are the same,  
 391  $m_1 = \dots = m_n = m$ . In all three settings, 200 datasets were generated for each combina-  
 392 tion of  $n$  and  $m$ . We focused on settings where  $m$  is small compared to  $n$ , to test the scope  
 393 of the theory in Section 4. For all simulations, the power parameter was fixed at  $\kappa = 2$ ,  
 394 while the hybrid estimator was obtained by refitting the selected submodel using adaptive  
 395 quadrature via the R package `lme4` (Bates et al., 2015).

396 For each setting, performance was assessed by the percentage of correctly chosen over-  
 397 all models, fixed effects, and random effects, as well as several measures of fit. Let  $\hat{\Psi}_{\text{method}}$   
 398 and  $\hat{\mathbf{b}}_{\text{method}}$  generically denote the parameter estimates and predicted random effects ob-  
 399 tained directly from regularized PQL or the hybrid estimation approach discussed in Sec-  
 400 tion 3.3. Then for both estimation methods we calculated the following four quantities:  
 401 mean absolute bias of the estimates  $\text{E} \left( \|\hat{\Psi}_{\text{method}} - \Psi_0\|_1 \right)$  where  $\|\cdot\|_1$  denotes the  $L_1$  norm,  
 402 total variance of the estimates  $\sum_{l=1}^{\dim(\Psi)} \text{Var}(\hat{\Psi}_{\text{method},l})$ , mean squared prediction error for ran-  
 403 dom effects  $\text{E}(\|\hat{\mathbf{b}}_{\text{method}} - \mathbf{b}_0\|^2)$ , and the mean predicted log-likelihood  $\text{E}\{\ell_{\text{pred}}(\hat{\Psi}_{\text{method}})\}$

404 evaluated using a validation dataset. For all four quantities, the expectations and variances  
405 were calculated empirically across the simulated datasets. Afterwards, for each quantity  
406 we constructed a ratio comparing the hybrid estimation approach to estimates directly from  
407 regularized PQL, such that ratios less than one imply the hybrid estimator has lower ab-  
408 solute bias/total variance/prediction error/predicted log-likelihood relative to regularized  
409 PQL.

## 410 5.1 Normal Responses

411 We replicated the design of Bondell et al. (2010), which was subsequently used by Fan  
412 and Li (2012) and Lin et al. (2013), so we can compare our method with other recently  
413 proposed penalized likelihood methods for linear mixed models. Datasets were generated  
414 based on the true model  $y_{ij} \sim N(\mathbf{x}_{ij}^T \boldsymbol{\beta} + \mathbf{z}_{ij}^T \mathbf{b}_i, \sigma^2)$ , where  $p_f = 9$  fixed effects with  
415 fixed intercept,  $p_r = 4$  random effects including a random intercept, and  $\sigma^2 = 1$ . The  
416 vector of true fixed effects parameters was set to  $\boldsymbol{\beta}_0 = (1, 1, 0, \dots, 0)$ , while the true  $4 \times 4$   
417 random effect covariance matrix is given by  $\text{vech}(\mathbf{D}_0) = (9, 4.8, 0.6, 0, 4, 1, 0, 1, 0, 0)$ . In  
418 other words, there were seven uninformative fixed effects and one uninformative random  
419 effect. All the elements of  $\mathbf{x}_{ij}$  and the last three elements  $\mathbf{z}_{ij}$  were generated from the  
420 uniform distribution  $U[-2, 2]$ , with the first element of  $\mathbf{z}_{ij}$  set equal to one. Four penalized  
421 likelihood methods were compared: 1) regularized PQL estimation (rPQL), 2) the SCAD-P  
422 approach of Fan and Li (2012) using the SCAD penalty, 3) the M-ALASSO approach of  
423 Bondell et al. (2010) using an adaptive lasso, and 4) the two-stage ALASSO approach of  
424 Lin et al. (2013). The results for methods 2 to 4 were taken from their respective papers.

425 Regularized PQL performed strongly overall; it was the best at selecting both the cor-  
426 rect overall model and fixed effects in the small sample case, while in the large sample case  
427 there was little difference between it and SCAD-P, which correctly identified the best model

Table 1: Results from simulation Setting 1 for linear mixed models. The methods are: regularized PQL (rPQL), SCAD-P (Fan and Li, 2012), M-ALASSO (Bondell et al., 2010), and ALASSO (Lin et al., 2013). Performance was assessed in terms of percentage datasets with correctly chosen overall models (%C), fixed effects (%CF), and random effects (%CR), as well as the ratios of mean absolute bias (Bias) and total variance (Var) of the estimates, mean squared prediction error (PSE), and predicted log-likelihood (PL).

$(n, m)$	Method	%C	% CF	% CR	Bias/Var/PSE/PL
(30, 5)	rPQL	88	98	88	0.84/1.02/0.88/0.97
	SCAD-P	-	90	86	-
	M-ALASSO	71	73	79	-
	ALASSO	79	81	96	-
(60, 10)	rPQL	98	99	98	0.99/1.03/0.97/0.95
	SCAD-P	100	100	100	-
	M-ALASSO	83	83	89	-
	ALASSO	95	96	99	-

428 in all simulated datasets (Table 1). The fact that the performance of regularized PQL was  
429 closer to SCAD-P than the other two penalized methods was not surprising, as regularized  
430 PQL and SCAD-P adopt a similar approach to group penalizing the random effect coeffi-  
431 cients, while M-ALASSO and ALASSO instead penalize the Cholesky decomposition of  
432 the random effect covariance matrix.

433 All the ratios were relatively close to one, suggesting that there was no substantial dif-  
434 ferences between the hybrid estimation approach compared to regularized PQL. This was  
435 not surprising, given that for linear mixed models, PQL estimation does produce asymp-  
436 totically unbiased and consistent estimators even in the setting where  $m$  is fixed (Breslow  
437 and Clayton, 1993). On computation time, regularized PQL took an average of 26 and 59  
438 seconds to fit the  $(n = 30, m = 5)$  and  $(n = 60, m = 10)$  settings respectively. We believe  
439 these times are competitive, while acknowledging that further reductions could have been  
440 made if we had used more sophisticated methods of optimization.

## 5.2 Bernoulli Responses

We simulated datasets from a Bernoulli GLMM using a logit link, with  $p_f = p_r = 9$  covariates, both including an intercept term. For  $i = 1, \dots, n$ , vectors of fixed effect covariates  $\mathbf{x}_{ij}$  were constructed with a one in the first term and the remaining terms generated from a multivariate normal distribution  $N_8(\mathbf{0}, \Sigma)$  with  $\Sigma_{rs} = 0.5^{|r-s|}$ . The random effect covariates  $\mathbf{z}_{ij}$  were set equal to  $\mathbf{x}_{ij}$ . The vector of true fixed effects parameters was set to  $\beta_0 = (-0.1, 1, -1, 1, -1, 0, \dots, 0)$ , while the true random effect covariance matrix was a  $9 \times 9$  diagonal matrix with the first three diagonal elements set to  $(3, 2, 1)$  and the remaining diagonal entries zero.

We are not aware of any available software for penalized joint selection in GLMMs. For comparison then, we write our own code to implement the following two penalized likelihood methods: 1) extending the M-ALASSO penalty of Bondell et al. (2010) to the case of non-Gaussian responses, with the tuning parameter chosen using their recommended BIC, 2) the adaptive lasso penalty of Ibrahim et al. (2011), with the tuning parameters chosen using their proposed  $IC_Q$  criterion. Estimation for both methods was performed using a penalized Monte-Carlo EM algorithm and, due to their heavy computational load, considered a sequence (grid) of 100 values (combinations) of the tuning parameter. Aside from these two penalties, we also applied the `glmLasso` package (Groll and Tutz, 2014), which performs fixed effects selection *only* in GLMMs using the unweighted lasso penalty. Since `glmLasso` only performs fixed effects selection, we assumed that the random effects component was known, i.e. only the first three elements of  $\mathbf{z}_{ij}$  were included in the random effects structure. As recommended by Groll and Tutz (2014), BIC was used to select the tuning parameter in `glmLasso`.

Finally, as an alternative to penalized likelihood, we included for comparison a two stage, forward selection method using  $BIC(\alpha) = -2\ell(\tilde{\Psi}_\alpha) + \log(N) \dim(\tilde{\Psi}_\alpha)$ , where  $\tilde{\Psi}_\alpha$

466 denotes the maximum likelihood estimates for for submodel  $\alpha$ . At the first stage, a saturated  
467 fixed effects structure was assumed and forward selection performed on the random effects.  
468 At the second stage, all random effects chosen in the first stage were entered into the model  
469 as fixed effects also, and forward selection was used on the remaining covariates to select  
470 them as fixed effects only. Compared to all subsets selection, the two stage approach is  
471 not only computationally more efficient, but also preserves the hierarchy of the covariates  
472 present in longitudinal GLMMs (Hui et al., 2016).

473 Regularized PQL performed best at selecting both fixed and random effects, with per-  
474 formance improving with  $m$  and  $n$  (Table 2). Comparing the hybrid and regularized PQL  
475 estimation methods, we see that the hybrid estimator produces considerably less biased but  
476 more variable estimates. This is consistent with the effects of penalization, that is, shrink-  
477 age of the fixed and random effects will reduce the variability of the estimates at the expense  
478 of increased bias. On the other hand, both the ratios for mean squared error prediction and  
479 predictive log-likelihood are less than one, particularly when  $m$  is small compared to  $n$ ,  
480 suggesting that the hybrid estimator did have improved predictive performance compared  
481 to directly using the regularized PQL estimates. The M-ALASSO penalty, `glmmlasso`,  
482 and forward selection using BIC all performed slightly poorer than regularized PQL at  
483 selecting the fixed effects, while on random effects selection M-ALASSO and forward se-  
484 lection using BIC had a tendency to overfit. Finally, the penalty of Ibrahim et al. (2011)  
485 performed poorly in this simulation, with subsequent investigation revealing that  $IC_Q$  al-  
486 most always chose the smallest possible set of tuning parameters (leading to the saturated  
487 model being selected). It also tended to behave erratically e.g., the loss function compo-  
488 nent of  $IC_Q$  did not vary monotonically with model complexity. It should be noted that  $IC_Q$   
489 criterion was, in fact, *not* recommended for use by the authors in an earlier paper (Ibrahim  
490 et al., 2008).

Table 2: Results from simulation Setting 2 for Bernoulli GLMMs. The methods are: regularized PQL (rPQL), M-ALASSO (Bondell et al., 2010), I-ALASSO (Ibrahim et al., 2011), `glmmLasso` (Groll and Tutz, 2014), and forward selection (Forward Sel.) using  $BIC(\alpha)$ . Performance was assessed in terms of the percentage of datasets with correctly chosen overall models (%C), fixed effects (%CF), and random effects (%CR), as well as ratios of mean absolute bias (Bias) and total variance (Var) of the estimates, mean squared prediction error (PSE), and predicted log-likelihood (PL). Finally, the mean computation time for each method was also recorded, with standard deviations in parentheses.

$(n, m)$	Method	%C	% CF	% CR	Comp. time	Bias/Var/PSE/PL
(50, 10)	rPQL	67	93	67	238(65)	0.21/18.28/0.73/0.71
	M-ALASSO	12	56	17	$\approx 10^4$	-
	I-ALASSO	0	0	0	7309(884)	-
	<code>glmmLasso</code>	-	73	-	908(109)	-
	Forward Sel.	9	94	10	192(67)	-
(50, 20)	rPQL	86	94	90	256(33)	0.27/4.70/0.63/0.85
	M-ALASSO	37	86	44	$\approx 10^4$	-
	I-ALASSO	0	10	0	8748(1115)	-
	<code>glmmLasso</code>	-	89	-	1301(162)	-
	Forward Sel.	74	98	76	1686(391)	-
(100, 10)	rPQL	78	96	81	390(73)	0.08/6.34/0.69/0.76
	M-ALASSO	12	83	17	$\approx 20^4$	-
	I-ALASSO	0	1	0	$\approx 10^4$	-
	<code>glmmLasso</code>	-	78	-	3226(231)	-
	Forward Sel.	33	93	36	1614(326)	-
(100, 20)	rPQL	95	97	98	501(98)	0.21/4.07/0.70/0.86
	M-ALASSO	43	92	45	$\approx 20^4$	-
	I-ALASSO	0	18	0	$\approx 10^4$	-
	<code>glmmLasso</code>	-	94	-	5738(264)	-
	Forward Sel.	95	97	98	6493(1031)	-

491 Except for  $(n = 50, m = 10)$  where forward selection using BIC was slightly faster,  
492 regularized PQL was also the fastest method at performing joint selection, with compu-  
493 tation time typically an order of magnitude smaller than the four competing approaches  
494 (Table 2). The long computation times of M-ALASSO and the penalty of Ibrahim et al.  
495 (2011) could be attributed to the need for a penalized Monte-Carlo EM algorithm, in con-  
496 trast to regularized PQL which does not involve any integration. Finally, computation times  
497 for forward selection using BIC scaled the worst with  $n$  and  $m$  e.g., doubling the cluster  
498 size from  $m = 10$  to 20 led to at least four-fold increase in estimation time.

499 Simulation results for the Poisson GLMMs are presented in the Supplementary Mate-  
500 rial, and present similar trends to those seen in the Bernoulli GLMM design above. That  
501 is, regularized PQL performed competitively in jointly selecting the fixed and random ef-  
502 fects, while taking much less time to fit the solution path than competing methods. Also  
503 presented in the Supplementary Material are results based on using forward selection with  
504 other types of information criteria, which performed worse than  $\text{BIC}(\alpha)$  shown above, as  
505 well as simulation designs where  $m$  explicitly grows as a function of  $n$ , which empirically  
506 confirmed the estimation and selection consistency established in Section 4.

## 507 **6 Application to Forest Health Monitoring**

508 We applied regularized PQL estimation to a longitudinal dataset on the health status of  
509 beech trees at plots located across northern Bavaria, Germany. The aim of the analysis was  
510 to uncover important baseline and time varying covariates influencing the probability of a  
511 tree experiencing defoliation.



Table 3: Nine baseline (time independent) and two time varying covariates available for selection in the forest health dataset.

Covariates	Brief description
<i>Baseline covariates</i>	
Alkali	Proportion of base alkali ions; categorical (very low, low, high, very high)
Canopy	Forest canopy density; continuous (%)
Elevation	Elevation above sea level; continuous (meters)
Fertilization	Fertilization applied; binary (yes, no)
Humus	Humus layer thickness; ordinal (five levels)
Inclination	Slope inclination; continuous (%)
Moisture	Soil moisture level; categorical (moderately dry, moderately moist, moist)
Soil	Soil layer depth; continuous (centimeters)
Stand	Stand type; categorical (deciduous, mixed)
<i>Time varying covariates</i>	
Age	Age of observation stands; continuous (years)
pH	Soil pH at 0–2 centimeters; continuous (centimeters)

512 Different versions of the data, i.e. with differing predictors and response type, have  
513 been considered previously by Kneib et al. (2009), who focused on the spatial effects, and  
514 Groll and Tutz (2014), who examined this data to illustrate high-dimensional GLMMs. In  
515 particular, Groll and Tutz (2014) also had the goal of identifying important predictors of  
516 tree defoliation, and we will compare our results with theirs. The version of the dataset we  
517 used is the `ForestHealth` dataset in the `R2BayesX` package (Belitz et al., 2015). The  
518 dataset consists of  $n = 78$  trees with  $m = 22$  measurements for all trees, with a binary  
519 response  $y_{ij} = 1$  indicating that defoliation exceeding 12.5% and  $y_{ij} = 0$  otherwise. As  
520 displayed in Table 3, nine baseline and two time varying covariates were recorded. All  
521 continuous covariates were standardized prior to analysis, while dummy variables were  
522 created for the categorical variables.

523 We fitted a Bernoulli GLMM with all covariates included as fixed effects. Furthermore,  
524 to account for any potential non-linear relationship between age and the probability of

525 defoliation on the logit scale, we included polynomial terms for age as fixed effects up to  
 526 the fourth power, similar to Groll and Tutz (2014). For the random effects, we included a  
 527 random intercept to account for heterogeneity in the overall health of the trees at baseline,  
 528 and random slopes for age and pH to capture the variability between trees in their response  
 529 to these covariates over time. We chose not to include any polynomial terms as random  
 530 effects for ease of interpretation. We first fitted a saturated model to construct adaptive  
 531 lasso weights. Then we used regularized PQL with the  $IC(\lambda)$  in (3) to perform model  
 532 selection, where  $IC(\lambda)$  was used to select both  $\lambda$  and  $\kappa$ , the latter chosen from the range  
 533  $\{1, 2\}$ . This resulted in the model

$$\begin{aligned} \text{logit}(\mu_{ij}) &= 0.528 + 0.364\text{Age}_{ij} - 1.235\text{Canopy}_i - 0.101\text{pH}_{ij} \\ &\quad + b_{0i} + b_{1i}\text{Age}_{ij} + b_{2i}\text{pH}_{ij}; \quad i = 1, \dots, 78, j = 1, \dots, 22 \\ \widehat{\text{Cov}}(\mathbf{b}_i) &= \begin{pmatrix} 5.042 & 2.822 & 1.024 \\ - & 3.427 & 0.928 \\ - & - & 0.839 \end{pmatrix}. \end{aligned}$$

534 Not surprisingly, older trees, increased soil acidity (lower pH), and denser forest canopy  
 535 cover were all associated with increased probability of defoliation. There was substantial  
 536 heterogeneity in the baseline status of the trees (remembering the continuous covariates  
 537 were standardized), as well as in their responses to age and pH. Regularized PQL shrunk  
 538 all the polynomial terms of age to zero, suggesting that perhaps any truly non-linear effect  
 539 of age was masked by the large variability between trees in their linear responses and/or that  
 540 the non-linear effects were comparatively small compared to the between-tree variability.  
 541 To confirm this, we fitted the selected submodel in the R package `lme4` using Laplace's  
 542 approximation, and compared it to a GLMM that included polynomial terms for age up the

543 fourth power. The resulting bootstrap likelihood ratio test confirmed that these polynomial  
544 fixed effects for age were indeed not significant (P-value = 0.11). Finally, all the off-  
545 diagonal terms in the estimated random effect covariance matrix were positive, indicating  
546 that large effects for one predictor tended to occur with large effects in the other predictors,  
547 e.g., the higher the baseline probability of defoliation, the worse the effect of increasing  
548 age and soil acidity on the the tree's health.

549 The results obtained here differ from those in Groll and Tutz (2014), who applied the  
550 `glmLasso` package to a very similar version of this dataset, in some important ways:  
551 1) Groll and Tutz (2014) did not have pH as a predictor in their analysis, whereas we  
552 found that, based on regularized PQL, pH was both an important fixed and random effect;  
553 2) the method of Groll and Tutz (2014) identified an important fixed, quadratic effect of  
554 age, although the magnitude of the coefficient was very close to zero; 3) regularized PQL  
555 identified canopy cover as an important predictor, whereas Groll and Tutz (2014) identi-  
556 fied stand type as important. Perhaps the driving reason behind these differences was that  
557 `glmLasso` selects only fixed effects, and Groll and Tutz (2014) only included a random  
558 intercept in the model. By contrast, regularized PQL performs joint selection so we could  
559 include and select on Age and pH as random slopes, and indeed both these covariates were  
560 identified as being significant effects.

## 561 **7 Discussion**

562 Joint selection of fixed and random effects in mixed models is a challenging problem, due to  
563 both the intractability of the marginal likelihood and the large number of candidate models.  
564 In this article, we proposed regularized PQL estimation to overcome these problems. By  
565 combining the PQL with adaptive lasso penalties for selecting the fixed and random effects,

566 regularized PQL offers a attractive method of computing the solution path. We showed  
567 that regularized PQL is selection consistent. This is an important result given PQL was  
568 originally motivated by Breslow and Clayton (1993) as a fast but approximate method of  
569 estimating GLMMs. With regularized PQL, we have a computationally fast approach of  
570 joint selection that asymptotically selects the true set of fixed and random effects. We  
571 proposed an information criterion for choosing the tuning parameter in regularized PQL  
572 which leads to selection consistency. The criterion combines a BIC-type model complexity  
573 penalty for the fixed effects with a AIC-type penalty for the random effects. This is a  
574 reflection of the differing degrees of model complexity needed for the fixed coefficients,  
575 which grows at rate  $O(1)$ , versus the random coefficients, which grows at rate  $O(n)$ . In the  
576 linear regression and penalized GLM contexts respectively, Shao (1997) and Zhang et al.  
577 (2010) investigated the impacts of differing degrees of model complexity, and our criterion  
578 can be regarded as an extension of such results to the GLMM context using regularized  
579 PQL estimation.

580 Simulations demonstrate the selection consistency of regularized PQL in conjunction  
581 with the proposed information criterion, showing that it can outperform other methods of  
582 joint selection while offering considerable reductions in computation time. The use of a  
583 hybrid estimation method further helps to reduce finite sample bias and improve predic-  
584 tion. Indeed, using regularized PQL for fast model selection only mirrors other works in  
585 the GLM context, where penalized likelihood approaches have been proposed purely as  
586 a means of computationally efficient model selection (e.g., Schelldorfer et al., 2014; Hui  
587 et al., 2015). Of course, we acknowledge further simulations are required to fully assess  
588 the robustness of regularized PQL selection e.g., how it performs when the truly non-zero  
589 coefficients and hence signal to noise ratio is small, and that there are other methods of joint  
590 selection in GLMMs which were not included in our study e.g., the predictive shrinkage

591 selection method of Hu et al. (2015) designed specifically for Poisson mixture models in  
592 the context of network analysis.

593 One obvious extension to make to regularized PQL estimation is to high-dimensional  
594 GLMMs, where the number of fixed and random effects grows with the number of clusters  
595 and/or cluster size; see for example the recent works of Fan and Li (2012) and Groll and  
596 Tutz (2014). For the case where  $p_r$  remains fixed but  $p_f$  is permitted to grow, we believe  
597 the estimation and selection consistency results established in this article will continue to  
598 hold, provided the conditions on the tuning parameter are altered slightly. In a more general  
599 setting where  $p_r$  grows with  $m_1$  and  $n$ , some of the results established for high-dimensional  
600 penalized GLMs (see the overview by Fan and Lv, 2010) may in principle be adapted to  
601 GLMMs, especially since the PQL treats the random effects as if they are fixed coefficients.  
602 Another possible extension which is especially useful for longitudinal studies is to modify  
603 rPQL so that the penalties select covariates in a hierarchical manner, such that all covariates  
604 in the model are chosen as either fixed effects only or composite (fixed and random) effects  
605 (see for instance, Hui et al., 2016). This reflects the notion that covariates in longitudinal  
606 GLMMs should not be included in the model as random slopes only.

## 607 **Acknowledgements**

608 This research was supported by the Australian Research Council discovery project grant  
609 DP140101259. Thanks to Andreas Groll and Gordana Popovic for useful discussions.

## 610 **References**

- 611 Bates, D., Maechler, M., Bolker, B., and Walker, S. (2015). *lme4: Linear mixed-effects*  
612 *models using Eigen and S4*. R package version 1.1-8.
- 613 Belitz, C., Brezger, A., Kneib, T., Lang, S., and Umlauf, N. (2015). *BayesX: Software for*  
614 *Bayesian inference in structured additive regression models*. Version 1.0.
- 615 Bondell, H. D., Krishna, A., and Ghosh, S. K. (2010). Joint variable selection for fixed and  
616 random effects in linear mixed-effects models. *Biometrics*, 66:1069–1077.
- 617 Breslow, N. E. and Clayton, D. G. (1993). Approximate inference in generalized linear  
618 mixed models. *Journal of the American Statistical Association*, 88:9–25.
- 619 Demidenko, E. (2013). *Mixed Models: Theory and Applications with R*. Wiley Series in  
620 Probability and Statistics. Wiley.
- 621 Fan, J. and Li, R. (2001). Variable selection via nonconcave penalized likelihood and its  
622 oracle properties. *Journal of the American Statistical Association*, 96:1348–1360.
- 623 Fan, J. and Lv, J. (2010). A selective overview of variable selection in high dimensional  
624 feature space. *Statistica Sinica*, 20:101–148.
- 625 Fan, Y. and Li, R. (2012). Variable selection in linear mixed effects models. *The Annals of*  
626 *Statistics*, 40:2043–2068.
- 627 Friedman, J., Hastie, T., and Tibshirani, R. (2010). Regularization paths for generalized  
628 linear models via coordinate descent. *Journal of Statistical Software*, 33:1–22.
- 629 Garcia, T. P., Müller, S., Carroll, R. J., and Walzem, R. L. (2014). Identification of impor-

630 tant regressor groups, subgroups and individuals via regularization methods: application  
631 to gut microbiome data. *Bioinformatics*, 30:831–837.

632 Groll, A. and Tutz, G. (2014). Variable selection for generalized linear mixed models by  
633  $\ell_1$ -penalized estimation. *Statistics and Computing*, 24:137–154.

634 Hu, K., Choi, J., Sim, A., and Jiang, J. (2015). Best predictive generalized linear mixed  
635 model with predictive lasso for high-speed network data analysis. *International Journal*  
636 *of Statistics and Probability*, 4:132–148.

637 Hui, F. K. C., Mueller, S., and Welsh, A. H. (2016). Hierarchical Selection of Fixed and  
638 Random Effects in Generalized Linear Mixed Models. *Statistica Sinica*, Accepted for  
639 publication.

640 Hui, F. K. C., Warton, D. I., and Foster, S. D. (2015). Tuning parameter selection for the  
641 adaptive lasso using ERIC. *Journal of the American Statistical Association*, 110:262–  
642 269.

643 Ibrahim, J. G., Zhu, H., Garcia, R. I., and Guo, R. (2011). Fixed and random effects  
644 selection in mixed effects models. *Biometrics*, 67:495–503.

645 Ibrahim, J. G., Zhu, H., and Tang, N. (2008). Model selection criteria for missing-data  
646 problems using the EM algorithm. *Journal of the American Statistical Association*,  
647 103:1648–1658.

648 Jiang, J., Jia, H., and Chen, H. (2001). Maximum posterior estimation of random effects in  
649 generalized linear mixed models. *Statistica Sinica*, 11:97–120.

650 Jiang, J., Rao, J. S., Gu, Z., Nguyen, T., et al. (2008). Fence methods for mixed model  
651 selection. *The Annals of Statistics*, 36:1669–1692.

- 652 Kneib, T., Hothorn, T., and Tutz, G. (2009). Variable selection and model choice in geoad-  
653 ditive regression models. *Biometrics*, 65:626–634.
- 654 Lehmann, E. (1983). *Theory of Point Estimation*. Wiley.
- 655 Lele, S. R., Nadeem, K., and Schmuland, B. (2010). Estimability and likelihood infer-  
656 ence for generalized linear mixed models using data cloning. *Journal of the American*  
657 *Statistical Association*, 105:1617–1625.
- 658 Lin, B., Pang, Z., and Jiang, J. (2013). Fixed and random effects selection by REML and  
659 pathwise coordinate optimization. *Journal of Computational and Graphical Statistics*,  
660 22:341–355.
- 661 Lin, X. and Breslow, N. E. (1996). Bias correction in generalized linear mixed models with  
662 multiple components of dispersion. *Journal of the American Statistical Association*,  
663 91:1007–1016.
- 664 McCulloch, C. E. (1997). Maximum likelihood algorithms for generalized linear mixed  
665 models. *Journal of the American Statistical Association*, 92:162–170.
- 666 Müller, S., Scaely, J. L., and Welsh, A. H. (2013). Model selection in linear mixed models.  
667 *Statistical Science*, 28:135–167.
- 668 Müller, S. and Welsh, A. (2009). Robust model selection in generalized linear models.  
669 *Statistica Sinica*, 19:1155–1170.
- 670 Pan, J. and Huang, C. (2014). Random effects selection in generalized linear mixed models  
671 via shrinkage penalty function. *Statistics and Computing*, 24:725–738.
- 672 Peng, H. and Lu, Y. (2012). Model selection in linear mixed effect models. *Journal of*  
673 *Multivariate Analysis*, 109:109–129.



- 674 Rabe-Hesketh, S., Skrondal, A., and Pickles, A. (2002). Reliable estimation of generalized  
675 linear mixed models using adaptive quadrature. *The Stata Journal*, 2:1–21.
- 676 Schelldorfer, J., Meier, L., and Bühlmann, P. (2014). Glmmlasso: an algorithm for high-  
677 dimensional generalized linear mixed models using  $\ell_1$ -penalization. *Journal of Compu-  
678 tational and Graphical Statistics*, 23:460–477.
- 679 Shao, J. (1997). An asymptotic theory for linear model selection. *Statistica Sinica*, 7:221–  
680 264.
- 681 Tierney, L. and Kadane, J. B. (1986). Accurate approximations for posterior moments and  
682 marginal densities. *Journal of the American Statistical Association*, 81:82–86.
- 683 Vaida, F. and Blanchard, S. (2005). Conditional Akaike information for mixed-effects  
684 models. *Biometrika*, 92:351–370.
- 685 Vonesh, E. F., Wang, H., Nie, L., and Majumdar, D. (2002). Conditional second-order gen-  
686 eralized estimating equations for generalized linear and nonlinear mixed-effects models.  
687 *Journal of the American Statistical Association*, 97:271–283.
- 688 White, H. (1982). Maximum Likelihood Estimation of Misspecified Models. *Economet-  
689 rica*, 50:1–26.
- 690 Zhang, Y., Li, R., and Tsai, C. (2010). Regularization parameter selections via generalized  
691 information criterion. *Journal of the American Statistical Association*, 105:312–323.

## 692 **A Derivation of Covariance Matrix update in equation (2)**

693 Consider the Laplace approximation log-likelihood  $\ell_{\text{LA}}(\Psi)$  given in Section 2 of the main  
 694 text. Substituting in the regularized PQL estimates  $\hat{\beta}_\lambda$  and  $\hat{\mathbf{b}}_\lambda$ , we obtain

$$\begin{aligned} \ell(\mathbf{D}) &= -\frac{n}{2} \log \det(\mathbf{D}) - \frac{1}{2} \sum_{i=1}^n \log \det(\mathbf{Z}_i^T \hat{\mathbf{W}}_{\lambda_i} \mathbf{Z}_i + \mathbf{D}^-) + \sum_{i=1}^n \sum_{j=1}^{m_i} \log f(y_{ij} | \hat{\beta}_\lambda, \hat{\mathbf{b}}_{\lambda_i}) \\ &\quad - \frac{1}{2} \sum_{i=1}^n \hat{\mathbf{b}}_{\lambda_i}^T \mathbf{D}^- \hat{\mathbf{b}}_{\lambda_i}, \end{aligned}$$

695 where for  $i = 1, \dots, n$ ,  $\hat{\mathbf{b}}_{\lambda_i}$  are the regularized PQL estimates of the random effects and  
 696  $\hat{\mathbf{W}}_{\lambda_i, jj} = (\text{Var}(y_{ij}) g'(\hat{\mu}_{\lambda_{ij}})^2)^-$ . Differentiating  $\ell(\mathbf{D})$  with respect to  $\mathbf{D}$ , we have

$$\begin{aligned} \frac{\partial \ell(\mathbf{D})}{\partial \text{vech}(\mathbf{D})} &= -\frac{n}{2} \text{vech}(\mathbf{D}^-) + \frac{1}{2} \sum_{i=1}^n (\mathbf{D}^- \otimes \mathbf{D}^-) \text{vech}\{(\mathbf{Z}_i^T \hat{\mathbf{W}}_{\lambda_i} \mathbf{Z}_i + \mathbf{D}^-)^{-1}\} \\ &\quad + \frac{1}{2} \sum_{i=1}^n (\mathbf{D}^- \otimes \mathbf{D}^-) \text{vech}(\hat{\mathbf{b}}_{\lambda_i} \hat{\mathbf{b}}_{\lambda_i}^T) \\ &= \mathbf{0} \end{aligned}$$

697 It follows that  $n(\mathbf{D}^- \otimes \mathbf{D}^-)^- \text{vech}(\mathbf{D}^-) = n \text{vech}(\mathbf{D}) = \sum_{i=1}^n \text{vech}\{(\mathbf{Z}_i^T \hat{\mathbf{W}}_{\lambda_i} \mathbf{Z}_i + \mathbf{D}^-)^{-1} +$   
 698  $\hat{\mathbf{b}}_{\lambda_i} \hat{\mathbf{b}}_{\lambda_i}^T\}$ , from which the formula in (2) of the main text follows.  $\square$

## 699 **B Outlines of Proofs**

700 Full derivations are found in the Supplementary Material; here we provide an outline for  
 701 each of these proofs.

702 Proof of Theorem 1: We consider the objective function  $\ell_p(\beta, \Gamma, \mathbf{b}) = \ell_{\text{PQL}}(\beta, \Gamma, \mathbf{b}) -$   
 703  $\lambda \sum_{k=1}^{p_f} v_k |\beta_k| - \lambda \sum_{l=1}^{p_r} w_l \|\mathbf{b}_{\bullet l}\|$  and define

704  $\Delta = n^{-1} \{ \ell_p(\boldsymbol{\beta}_0 + \alpha_m \mathbf{u}_1, \boldsymbol{\Gamma}, \mathbf{b}_0 + \alpha_m \mathbf{u}_2) - \ell_p(\boldsymbol{\beta}_0, \boldsymbol{\Gamma}, \mathbf{b}_0) \}$  for a vector  $\mathbf{u}$  of appropriate  
705 length and  $\alpha_m = m_1^{-1/2}$ . Under conditions (C1)-(C2), (C4) and (C5a), we show that this  
706 difference is asymptotically dominated by a quadratic term of form  $-(\alpha_m^2/2) \mathbf{u}^T \{ -n^{-1} \nabla^2 \ell_1(\boldsymbol{\beta}, \mathbf{b}) \} \mathbf{u}$ ,  
707 which is negative. This implies that with probability tending to one there exists a local max-  
708 imum at  $(\boldsymbol{\beta}_0, \mathbf{b}_0)$ , which we then show to be a global maximum.

709 Given the  $m_1^{1/2}$ -consistency from the first part of the theorem, to prove selection consis-  
710 tency of the regularized PQL estimates we need only show that for truly zero fixed and ran-  
711 dom effects, the signs of the score equations  $\partial \ell_p(\boldsymbol{\Psi}, \mathbf{b}) / \partial \beta_k |_{\hat{\boldsymbol{\Psi}}_\lambda, \hat{\mathbf{b}}_\lambda}$  and  $\partial \ell_p(\boldsymbol{\Psi}, \mathbf{b}) / \partial b_{il} |_{\hat{\boldsymbol{\Psi}}_\lambda, \hat{\mathbf{b}}_\lambda}$   
712 depend asymptotically only on the sign of the estimated coefficients. This is proved by ex-  
713 panding the score equations about the true parameter values and, in particular, using con-  
714 dition (C5b) to show that the derivative of the adaptive (group) lasso penalty dominates all  
715 the terms in the score equations.

716 Proof of Lemma 1: We consider separately the cases of underfitted and overfitted mod-  
717 els. In the first case, we utilize condition (C6) to show that the difference in the loss function  
718  $-2 \sum_{i=1}^n \sum_{j=1}^{m_i} \log f(y_{ij} | \tilde{\boldsymbol{\beta}}_\alpha, \tilde{\mathbf{b}}_{\alpha i})$  between any underfitted model and the true model is positive  
719 and asymptotically dominates all differences in the model complexity. In the second case,  
720 Condition (C3) is utilized to show that the difference in the loss function between any over-  
721 fitted model and the true model is asymptotically dominated by the difference in the model  
722 complexity penalties  $\log(n) \dim(\tilde{\boldsymbol{\beta}}_\alpha) + 2n \dim(\tilde{\mathbf{b}}_{\alpha i})$ , which by definition is greater than  
723 zero when overfitting.

724 Proof of Theorem 2: Under conditions (C1)-(C2), (C4), we prove the result  $N^{-1} \ell_{\text{PQL}}(\hat{\boldsymbol{\Psi}}_{\lambda_0}, \hat{\mathbf{b}}_{\lambda_0}) =$   
725  $N^{-1} \ell_1(\tilde{\boldsymbol{\beta}}_{\alpha_0}, \tilde{\mathbf{b}}_{\alpha_0}) + o_p(1)$ . We then show for any tuning parameter  $\lambda$  producing an under-  
726 fitted or overfitted model  $\alpha$ , it holds that  $\{ \text{IC}(\lambda) - \text{IC}(\lambda_0) \} \geq \{ \text{IC}_{\text{proxy}}(\alpha) - \text{IC}_{\text{proxy}}(\alpha_0) \}$ .  
727 Since the right hand side is positive with probability tending to one by Lemma 1, the result  
728 follows.