

Linear mixed models and penalized least squares

Douglas M. Bates

Department of Statistics, University of Wisconsin–Madison

Saikat DebRoy

Department of Biostatistics, Harvard School of Public Health

Abstract

Linear mixed-effects models are an important class of statistical models that are not only used directly in many fields of applications but also used as iterative steps in fitting other types of mixed-effects models, such as generalized linear mixed models. The parameters in these models are typically estimated by maximum likelihood (ML) or restricted maximum likelihood (REML). In general there is no closed form solution for these estimates and they must be determined by iterative algorithms such as EM iterations or general nonlinear optimization. Many of the intermediate calculations for such iterations have been expressed as generalized least squares problems. We show that an alternative representation as a penalized least squares problem has many advantageous computational properties including the ability to evaluate explicitly a profiled log-likelihood or log-restricted likelihood, the gradient and Hessian of this profiled objective, and an ECME update to refine this objective.

Key words: REML, gradient, Hessian, EM algorithm, ECME algorithm, maximum likelihood, profile likelihood, multilevel models

1 Introduction

We will establish some results for the penalized least squares representation of a general form of a linear mixed-effects model, then show how these results specialize for particular models. The general form we consider is

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{b} + \boldsymbol{\epsilon} \quad \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \sigma^2\mathbf{I}), \mathbf{b} \sim \mathcal{N}(\mathbf{0}, \sigma^2\boldsymbol{\Omega}^{-1}), \boldsymbol{\epsilon} \perp \mathbf{b} \quad (1)$$

where \mathbf{y} is the n -dimensional response vector, \mathbf{X} is an $n \times p$ model matrix for the p -dimensional fixed-effects vector $\boldsymbol{\beta}$, \mathbf{Z} is the $n \times q$ model matrix for

the q -dimensional random-effects vector \mathbf{b} that has a Gaussian distribution with mean $\mathbf{0}$ and relative precision matrix $\mathbf{\Omega}$ (i.e., $\mathbf{\Omega}$ is the precision of \mathbf{b} relative to the precision of $\boldsymbol{\epsilon}$), and $\boldsymbol{\epsilon}$ is the random noise assumed to have a spherical Gaussian distribution. The symbol \perp indicates independence of random variables. We assume that \mathbf{X} has full column rank and that $\mathbf{\Omega}$ is positive definite. (If \mathbf{X} is rank deficient or if $\mathbf{\Omega}$ is singular then the model can be transformed to an alternative model that fulfills the desired conditions.)

A relative precision factor, $\mathbf{\Delta}$, is any $q \times q$ matrix that satisfies $\mathbf{\Omega} = \mathbf{\Delta}^\top \mathbf{\Delta}$. One possible $\mathbf{\Delta}$ is the Cholesky factor of $\mathbf{\Omega}$ but others can be used. Because $\mathbf{\Omega}$ is positive definite, any $\mathbf{\Delta}$ will be non-singular. In general $\mathbf{\Delta}$ (and hence $\mathbf{\Omega}$) depend on a k -dimensional parameter vector $\boldsymbol{\theta}$. Typically \mathbf{Z} , $\mathbf{\Omega}$, and $\mathbf{\Delta}$ are very large and sparse (mostly zeros) while k , the dimension of $\boldsymbol{\theta}$, is small.

The likelihood for the parameters $\boldsymbol{\beta}$, $\boldsymbol{\theta}$, and σ^2 , given \mathbf{y} , in model (1) is

$$L(\boldsymbol{\beta}, \boldsymbol{\theta}, \sigma^2 | \mathbf{y}) = \int \frac{\sqrt{|\mathbf{\Omega}|}}{(2\pi\sigma^2)^{(n+q)/2}} \exp\left(\frac{\|\mathbf{y} - \mathbf{X}\boldsymbol{\beta} - \mathbf{Z}\mathbf{b}\|^2 + \mathbf{b}^\top \mathbf{\Omega} \mathbf{b}}{-2\sigma^2}\right) d\mathbf{b}. \quad (2)$$

The restricted (or residual) maximum likelihood (REML) estimates of $\boldsymbol{\theta}$ and σ^2 optimize a related criterion that can be written

$$L_R(\boldsymbol{\theta}, \sigma^2) = \int L(\boldsymbol{\beta}, \boldsymbol{\theta}, \sigma^2) d\boldsymbol{\beta}. \quad (3)$$

This integral is not the typical way to derive or to justify REML estimates but, as we shall see, it provides a convenient form in which to evaluate L_R .

In the next section we show that the integrals in (2) and (3) can be expressed succinctly using the solution to a penalized least squares problem. In particular, we derive a profiled log-likelihood and a profiled log-restricted likelihood that depend on $\boldsymbol{\theta}$ only. Maximizing the profiled log-likelihood (or profiled log-restricted-likelihood) is generally a much smaller and more stable optimization problem than attempting to optimize the log-likelihood for the full parameter vector. In section 3 we derive new expressions for the gradient and the Hessian of these profiled log-likelihoods. Using these derivatives, the profiled log-likelihood can be rapidly optimized, say by Newton steps, once the neighbourhood of $\hat{\boldsymbol{\theta}}$ has been determined.

Some heuristics provide reasonable starting estimates for $\boldsymbol{\theta}$ but these may not be sufficiently close to $\hat{\boldsymbol{\theta}}$ to ensure stable Newton steps. The expectation-maximization (EM) algorithm is a robust algorithm that approaches the neighbourhood of $\hat{\boldsymbol{\theta}}$ quickly but tends to converge slowly once it is in the neighbourhood. In section 4 we show how the penalized least squares results provide the update for a related algorithm called ‘‘expectation-conditional maximization-either’’ (ECME). This allows the starting estimates to be updated by a moderate number of ECME iterations before starting the Newton steps.

These results are derived for the general model (1). In sections 5 and 6 we consider several common special cases. Our implementation of some of these methods is discussed in section 7.

2 A Penalized Least-Squares Problem

For a fixed value of $\boldsymbol{\theta}$ we consider the penalized least squares problem defined by the augmented model matrix $\Phi(\boldsymbol{\theta})$ and the augmented response vector $\tilde{\mathbf{y}}$;

$$\min_{\mathbf{b}, \boldsymbol{\beta}} \left\| \tilde{\mathbf{y}} - \Phi(\boldsymbol{\theta}) \begin{bmatrix} \mathbf{b} \\ \boldsymbol{\beta} \end{bmatrix} \right\|^2 \quad \text{where } \Phi(\boldsymbol{\theta}) = \begin{bmatrix} \mathbf{Z} & \mathbf{X} \\ \Delta(\boldsymbol{\theta}) & \mathbf{0} \end{bmatrix} \quad \text{and } \tilde{\mathbf{y}} = \begin{bmatrix} \mathbf{y} \\ \mathbf{0} \end{bmatrix}. \quad (4)$$

One way to solve problem (4) is to form $\Phi_e = [\Phi, \tilde{\mathbf{y}}]$ and let \mathbf{R}_e be the Cholesky decomposition of $\Phi_e^T \Phi_e$

$$\Phi_e^T \Phi_e = \begin{bmatrix} \mathbf{Z}^T \mathbf{Z} + \boldsymbol{\Omega} & \mathbf{Z}^T \mathbf{X} & \mathbf{Z}^T \mathbf{y} \\ \mathbf{X}^T \mathbf{Z} & \mathbf{X}^T \mathbf{X} & \mathbf{X}^T \mathbf{y} \\ \mathbf{y}^T \mathbf{Z} & \mathbf{y}^T \mathbf{X} & \mathbf{y}^T \mathbf{y} \end{bmatrix} = \mathbf{R}_e^T \mathbf{R}_e \quad \text{where } \mathbf{R}_e = \begin{bmatrix} \mathbf{R}_{ZZ} & \mathbf{R}_{ZX} & \mathbf{r}_{Zy} \\ \mathbf{0} & \mathbf{R}_{XX} & \mathbf{r}_{Xy} \\ \mathbf{0} & \mathbf{0} & r_{yy} \end{bmatrix}. \quad (5)$$

The matrices \mathbf{R}_{ZZ} and \mathbf{R}_{XX} are upper triangular of dimension $q \times q$ and $p \times p$ respectively. The corresponding vectors, \mathbf{r}_{Zy} and \mathbf{r}_{Xy} , are of dimension q and p , and r_{yy} is a scalar. The conditions that $\boldsymbol{\Omega}$ be positive definite and \mathbf{X} have full column rank ensure that Φ has full column rank, and hence that \mathbf{R}_{ZZ} and \mathbf{R}_{XX} are nonsingular.

Representation (5) is a particular form of the mixed model equations described in Henderson (1984) [1]. We write the blocks in the opposite order from which they are typically written because of computational advantages associated with this order.

Using $\mathbf{a} = [-\mathbf{b}^T, -\boldsymbol{\beta}^T, 1]^T$ we can write the numerator of the exponent in the integral in (2) as

$$\begin{aligned} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta} - \mathbf{Z}\mathbf{b}\|^2 + \mathbf{b}^T \boldsymbol{\Omega} \mathbf{b} &= \mathbf{a}^T \Phi_e^T \Phi_e \mathbf{a} = \mathbf{a}^T \mathbf{R}_e^T \mathbf{R}_e \mathbf{a} = \|\mathbf{R}_e \mathbf{a}\|^2 \\ &= \|\mathbf{r}_{Zy} - \mathbf{R}_{ZX}\boldsymbol{\beta} - \mathbf{R}_{ZZ}\mathbf{b}\|^2 + \|\mathbf{r}_{Xy} - \mathbf{R}_{XX}\boldsymbol{\beta}\|^2 + r_{yy}^2 \end{aligned} \quad (6)$$

and a simple change of variable allows us to evaluate

$$\int \frac{1}{(2\pi\sigma^2)^{q/2}} \exp\left(\frac{\|\mathbf{r}_{Zy} - \mathbf{R}_{ZX}\boldsymbol{\beta} - \mathbf{R}_{ZZ}\mathbf{b}\|^2}{-2\sigma^2}\right) d\mathbf{b} = \frac{1}{\text{abs}|\mathbf{R}_{ZZ}|} = \frac{1}{\sqrt{|\mathbf{Z}^\top\mathbf{Z} + \boldsymbol{\Omega}|}} \quad (7)$$

Combining (2), (6), and (7) and taking the logarithm produces the log-likelihood $\ell(\boldsymbol{\beta}, \boldsymbol{\theta}, \sigma^2) = \log L(\boldsymbol{\beta}, \boldsymbol{\theta}, \sigma^2)$, which, for convenience, we write in the form of a deviance

$$-2\ell(\boldsymbol{\beta}, \boldsymbol{\theta}, \sigma^2) = \log\left(\frac{|\mathbf{Z}^\top\mathbf{Z} + \boldsymbol{\Omega}|}{|\boldsymbol{\Omega}|}\right) + n \log(2\pi\sigma^2) + \frac{r_{yy}^2 + \|\mathbf{r}_{Xy} - \mathbf{R}_{XX}\boldsymbol{\beta}\|^2}{\sigma^2} \quad (8)$$

leading to the following results for the maximum likelihood estimates (mle's):

- (1) The conditional mle of the fixed-effects, $\widehat{\boldsymbol{\beta}}(\boldsymbol{\theta})$, satisfies

$$\mathbf{R}_{XX}\widehat{\boldsymbol{\beta}}(\boldsymbol{\theta}) = \mathbf{r}_{Xy} \quad (9)$$

- (2) The conditional mle of the variance, $\widehat{\sigma}^2(\boldsymbol{\theta})$, is r_{yy}^2/n .
(3) The profiled log-likelihood, $\tilde{\ell}(\boldsymbol{\theta})$, a function of $\boldsymbol{\theta}$ only, is given by

$$\begin{aligned} -2\tilde{\ell}(\boldsymbol{\theta}) &= -2\ell(\widehat{\boldsymbol{\beta}}(\boldsymbol{\theta}), \boldsymbol{\theta}, \widehat{\sigma}^2(\boldsymbol{\theta})) \\ &= \log\left(\frac{|\mathbf{Z}^\top\mathbf{Z} + \boldsymbol{\Omega}|}{|\boldsymbol{\Omega}|}\right) + n \left[1 + \log\left(\frac{2\pi r_{yy}^2}{n}\right)\right] \end{aligned} \quad (10)$$

- (4) The conditional expected value of \mathbf{b} , which we write as $\widehat{\mathbf{b}}(\boldsymbol{\beta}, \boldsymbol{\theta})$, satisfies

$$\mathbf{R}_{ZZ}\widehat{\mathbf{b}}(\boldsymbol{\beta}, \boldsymbol{\theta}) = \mathbf{r}_{Zy} - \mathbf{R}_{ZX}\boldsymbol{\beta} \quad (11)$$

Typically we evaluate $\widehat{\mathbf{b}}(\widehat{\boldsymbol{\beta}}(\boldsymbol{\theta}), \boldsymbol{\theta})$, which we write as $\widehat{\mathbf{b}}(\boldsymbol{\theta})$.

- (5) The conditional distribution of \mathbf{b} is

$$\mathbf{b}|\mathbf{y}, \boldsymbol{\beta}, \boldsymbol{\theta}, \sigma^2 \sim \mathcal{N}\left(\widehat{\mathbf{b}}(\boldsymbol{\beta}, \boldsymbol{\theta}), \sigma^2 (\mathbf{Z}^\top\mathbf{Z} + \boldsymbol{\Omega})^{-1}\right) \quad (12)$$

2.1 REML results

As in (7), we can use a simple change of variable to obtain

$$\int \frac{1}{(2\pi\sigma^2)^{p/2}} \exp\left(\frac{\|\mathbf{r}_{Xy} - \mathbf{R}_{XX}\boldsymbol{\beta}\|^2}{-2\sigma^2}\right) d\boldsymbol{\beta} = \frac{1}{\text{abs}|\mathbf{R}_{XX}|} \quad (13)$$

providing the log-restricted-likelihood, $\ell_R(\boldsymbol{\theta}, \sigma^2) = \log L_R(\boldsymbol{\theta}, \sigma^2)$, as

$$-2\ell_R(\boldsymbol{\theta}, \sigma^2) = \log \left(\frac{|\mathbf{Z}^\top \mathbf{Z} + \boldsymbol{\Omega}| |\mathbf{R}_{XX}|^2}{|\boldsymbol{\Omega}|} \right) + (n-p) \log(2\pi\sigma^2) + \frac{r_{yy}^2}{\sigma^2}. \quad (14)$$

Noting that $|\mathbf{Z}^\top \mathbf{Z} + \boldsymbol{\Omega}| |\mathbf{R}_{XX}|^2 = |\boldsymbol{\Phi}^\top \boldsymbol{\Phi}|$ we have the following results:

- (1) The conditional REML estimate of the variance, $\widehat{\sigma}_R^2(\boldsymbol{\theta})$, is $r_{yy}^2/(n-p)$.
- (2) The profiled log-restricted-likelihood is given by

$$-2\tilde{\ell}_R(\boldsymbol{\theta}) = \log \left(\frac{|\boldsymbol{\Phi}^\top \boldsymbol{\Phi}|}{|\boldsymbol{\Omega}|} \right) + (n-p) \left[1 + \log \left(\frac{2\pi r_{yy}^2}{n-p} \right) \right]. \quad (15)$$

- (3) The conditional distribution of \mathbf{b} is

$$\mathbf{b} | \mathbf{y}, \boldsymbol{\theta}, \sigma^2 \sim \mathcal{N}(\widehat{\mathbf{b}}(\boldsymbol{\theta}), \sigma^2 \mathbf{V}_b) \quad (16)$$

where \mathbf{V}_b is the upper-left $q \times q$ submatrix of $(\boldsymbol{\Phi}^\top \boldsymbol{\Phi})^{-1}$,

$$\mathbf{V}_b = \mathbf{R}_{ZZ}^{-1} \left(\mathbf{I} + \mathbf{R}_{ZX} \mathbf{R}_{XX}^{-1} \mathbf{R}_{XX}^{-\top} \mathbf{R}_{ZX}^\top \right) \mathbf{R}_{ZZ}^{-\top} \quad (17)$$

2.2 A preliminary decomposition

Because we wish to evaluate the log-likelihood or log-restricted-likelihood for many different values of $\boldsymbol{\theta}$, we form a preliminary decomposition of the cross-products of the model matrices and the response,

$$\begin{bmatrix} \mathbf{Z}^\top \mathbf{Z} & \mathbf{Z}^\top \mathbf{X} & \mathbf{Z}^\top \mathbf{y} \\ \mathbf{X}^\top \mathbf{Z} & \mathbf{X}^\top \mathbf{X} & \mathbf{X}^\top \mathbf{y} \\ \mathbf{y}^\top \mathbf{Z} & \mathbf{y}^\top \mathbf{X} & \mathbf{y}^\top \mathbf{y} \end{bmatrix} = \mathbf{T}^\top \mathbf{T} \text{ where } \mathbf{T} = \begin{bmatrix} \mathbf{T}_{ZZ} & \mathbf{T}_{ZX} & \mathbf{t}_{Zy} \\ \mathbf{0} & \mathbf{T}_{XX} & \mathbf{t}_{Xy} \\ \mathbf{0} & \mathbf{0} & \mathbf{t}_{yy} \end{bmatrix} \quad (18)$$

so that the evaluation of (5) is equivalent to forming the orthogonal-triangular decomposition

$$\begin{bmatrix} \boldsymbol{\Delta} & \mathbf{0} & \mathbf{0} \\ \mathbf{T}_{ZZ} & \mathbf{T}_{ZX} & \mathbf{t}_{Zy} \\ \mathbf{0} & \mathbf{T}_{XX} & \mathbf{t}_{Xy} \\ \mathbf{0} & \mathbf{0} & \mathbf{t}_{yy} \end{bmatrix} = \mathbf{Q} \begin{bmatrix} \mathbf{R}_{ZZ} & \mathbf{R}_{ZX} & \mathbf{r}_{Zy} \\ \mathbf{0} & \mathbf{R}_{XX} & \mathbf{r}_{Xy} \\ \mathbf{0} & \mathbf{0} & r_{yy} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} \end{bmatrix}. \quad (19)$$

To ensure well-defined estimates for the parameters we should have $n \geq q + p$ which means that (18) produces a reduction, sometimes a substantial reduc-

tion, in the amount of data that must be stored and manipulated for each evaluation of the log-likelihood.

Furthermore, the matrices $\mathbf{Z}^\top \mathbf{Z}$ and \mathbf{T}_{ZZ} are sparse and we store and manipulate these matrices taking the sparsity into account (see §5 and §6 for specific examples). We will assume that $\mathbf{Z}^\top \mathbf{X}$, \mathbf{T}_{ZX} , $\mathbf{X}^\top \mathbf{X}$, and \mathbf{T}_{XX} are dense and are stored accordingly.

Equation (19) shows that t_{yy}^2 is a lower bound for r_{yy}^2 . Similarly, the residual sum of squares from regressing \mathbf{y} on \mathbf{X} is an upper bound on r_{yy}^2 . Thus the terms in (10 and 15) involving r_{yy}^2 are bounded.

The ratio $|\mathbf{Z}^\top \mathbf{Z} + \mathbf{\Omega}|/|\mathbf{\Omega}|$ is bounded below by unity (and approaches this bound as $\mathbf{\Omega}^{-1} \rightarrow \mathbf{0}$) so the profiled deviance (10) is bounded below by $n \left[1 + \log \left(2\pi t_{yy}^2/n \right) \right]$ and the profiled restricted deviance (15) is bounded below by $(n - p) \left[1 + \log \left(2\pi t_{yy}^2/(n - p) \right) \right]$.

As $\mathbf{\Omega}$ approaches singularity (say the minimum eigenvalue of $\mathbf{\Omega}$ approaches zero while the other eigenvalues are bounded above) $|\mathbf{Z}^\top \mathbf{Z} + \mathbf{\Omega}|/|\mathbf{\Omega}| \rightarrow \infty$. Thus we know that the ML or REML estimates will not occur on the boundary of the set of positive definite $\mathbf{\Omega}$.

It is possible that finite ML or REML estimates of $\mathbf{\Omega}$ will not exist. The minimum profiled deviance may correspond to an infinite precision (unbounded $\mathbf{\Omega}$), which is to say that the ML or REML estimates of $\sigma^2 \mathbf{\Omega}^{-1}$, the variance-covariance of \mathbf{b} , are singular.

3 Derivatives of the profiled log-likelihood

Expressions (10) and (15) provide extremely efficient ways to determine ML or REML estimates for a linear mixed-effects model because we can optimize these expressions as a function of $\boldsymbol{\theta}$ only, instead of as a function of the complete parameter vector $(\boldsymbol{\theta}^\top, \boldsymbol{\beta}^\top, \sigma^2)^\top$. The reduction in the dimension of the parameter over which we are optimizing helps to improve the performance and reliability of numerical optimization routines. Another way in which we can improve performance and reliability in numerical optimization is by providing analytic derivatives of the objective, which we do using results from Golub and Pereyra [2].

As in [2] we will use the operator \mathbf{D} to indicate the Frechet derivative of an array or a scalar. For example, $\mathbf{D} \Phi(\boldsymbol{\theta})$ is an array of dimension $n \times (q + p) \times k$, which we treat as k matrices of size $n \times (q + p)$ when writing expressions involving matrices and arrays. When we need to indicate partial derivatives

with respect to particular parameters we will use the notation $D_i \Phi(\boldsymbol{\theta})$ for the $n \times (q + p)$ matrix $\partial \Phi / \partial \theta_i$.

Golub and Pereyra [2] provide derivatives of the projection orthogonal to the column space of Φ , which, because $\Phi^T \Phi$ is nonsingular, we can write $P^\perp = \mathbf{I} - \Phi (\Phi^T \Phi)^{-1} \Phi^T$; the pseudo-inverse of Φ , which is $\Phi^+ = (\Phi^T \Phi)^{-1} \Phi^T$; and the residual sum of squares, $r_{yy}^2(\boldsymbol{\theta}) = \|P^\perp \tilde{\mathbf{y}}\|^2$. These derivatives are

$$D P^\perp = -P^\perp D \Phi \Phi^+ - (P^\perp D \Phi \Phi^+)^T, \quad (20)$$

$$\frac{1}{2} \nabla r_{yy}^2(\boldsymbol{\theta}) = -\tilde{\mathbf{y}}^T P^\perp D \Phi \Phi^+ \tilde{\mathbf{y}}, \quad (21)$$

$$D \Phi^+ = -\Phi^+ D \Phi \Phi^+ + \Phi^+ (\Phi^+)^T D \Phi^T P^\perp. \quad (22)$$

(Equation (22) is derived from equation (4.12) in [2], which has another term. However, that term is identically zero when Φ has full column rank.)

For \mathbf{A} a square, nonsingular matrix we have

$$\nabla(\log |\mathbf{A}(\boldsymbol{\theta})|) = \text{tr} [D(\mathbf{A}) \mathbf{A}^{-1}] \quad (23)$$

$$D(\mathbf{A}^{-1}) = -\mathbf{A}^{-1} D(\mathbf{A}) \mathbf{A}^{-1} \quad (24)$$

where tr denotes the trace of a matrix. Notice that (24) is a special case of (22).

In Appendix A we show that the gradient and Hessian terms for the penalized residual sum-of-squares, $r_{yy}^2(\boldsymbol{\theta})$, are

$$\nabla r_{yy}^2(\boldsymbol{\theta}) = \hat{\mathbf{b}}^T D \Omega \hat{\mathbf{b}} = \text{tr} (D \Omega \hat{\mathbf{b}} \hat{\mathbf{b}}^T) \quad (25)$$

$$D_j D_i r_{yy}^2(\boldsymbol{\theta}) = \hat{\mathbf{b}}^T [D_j (D_i \Omega) - 2D_j \Omega V_b D_i \Omega] \hat{\mathbf{b}} \quad (26)$$

providing the gradients

$$\nabla(-2\tilde{\ell}) = \text{tr} \left[D \Omega \left((\mathbf{Z}^T \mathbf{Z} + \Omega)^{-1} - \Omega^{-1} + \frac{\hat{\mathbf{b}} \hat{\mathbf{b}}^T}{\hat{\sigma} \hat{\sigma}} \right) \right] \quad (27)$$

$$\nabla(-2\tilde{\ell}_R) = \text{tr} \left[D \Omega \left(\mathbf{V}_b - \Omega^{-1} + \frac{\hat{\mathbf{b}} \hat{\mathbf{b}}^T}{\hat{\sigma}_R \hat{\sigma}_R} \right) \right] \quad (28)$$

and Hessian terms of the form

$$\begin{aligned} \mathbf{D}_j \mathbf{D}_i(-2\tilde{\ell}) &= \text{tr} \left[\mathbf{D}_j(\mathbf{D}_i \boldsymbol{\Omega}) \left((\mathbf{Z}^\top \mathbf{Z} + \boldsymbol{\Omega})^{-1} - \boldsymbol{\Omega}^{-1} + \frac{\hat{\mathbf{b}} \hat{\mathbf{b}}^\top}{\hat{\sigma} \hat{\sigma}} \right) \right] \\ &\quad - \text{tr} \left[\mathbf{D}_j \boldsymbol{\Omega} (\mathbf{Z}^\top \mathbf{Z} + \boldsymbol{\Omega})^{-1} \mathbf{D}_i \boldsymbol{\Omega} (\mathbf{Z}^\top \mathbf{Z} + \boldsymbol{\Omega})^{-1} \right] \end{aligned} \quad (29)$$

$$\begin{aligned} &\quad + \text{tr} \left(\mathbf{D}_j \boldsymbol{\Omega} \boldsymbol{\Omega}^{-1} \mathbf{D}_i \boldsymbol{\Omega} \boldsymbol{\Omega}^{-1} \right) - 2 \frac{\hat{\mathbf{b}}^\top}{\hat{\sigma}} \mathbf{D}_j \boldsymbol{\Omega} \mathbf{V}_b \mathbf{D}_i \boldsymbol{\Omega} \frac{\hat{\mathbf{b}}}{\hat{\sigma}} \\ &\quad - \frac{1}{n} \left(\frac{\hat{\mathbf{b}}^\top}{\hat{\sigma}} \mathbf{D}_j \boldsymbol{\Omega} \frac{\hat{\mathbf{b}}}{\hat{\sigma}} \right) \left(\frac{\hat{\mathbf{b}}^\top}{\hat{\sigma}} \mathbf{D}_i \boldsymbol{\Omega} \frac{\hat{\mathbf{b}}}{\hat{\sigma}} \right) \\ \mathbf{D}_j \mathbf{D}_i(-2\tilde{\ell}_R) &= \text{tr} \left[\mathbf{D}_j(\mathbf{D}_i \boldsymbol{\Omega}) \left(\mathbf{V}_b - \boldsymbol{\Omega}^{-1} + \frac{\hat{\mathbf{b}} \hat{\mathbf{b}}^\top}{\hat{\sigma}_R \hat{\sigma}_R} \right) \right] \\ &\quad - \text{tr} [\mathbf{D}_j \boldsymbol{\Omega} \mathbf{V}_b \mathbf{D}_i \boldsymbol{\Omega} \mathbf{V}_b] \\ &\quad + \text{tr} \left(\mathbf{D}_j \boldsymbol{\Omega} \boldsymbol{\Omega}^{-1} \mathbf{D}_i \boldsymbol{\Omega} \boldsymbol{\Omega}^{-1} \right) - 2 \frac{\hat{\mathbf{b}}^\top}{\hat{\sigma}_R} \mathbf{D}_j \boldsymbol{\Omega} \mathbf{V}_b \mathbf{D}_i \boldsymbol{\Omega} \frac{\hat{\mathbf{b}}}{\hat{\sigma}_R} \\ &\quad - \frac{1}{n-p} \left(\frac{\hat{\mathbf{b}}^\top}{\hat{\sigma}_R} \mathbf{D}_j \boldsymbol{\Omega} \frac{\hat{\mathbf{b}}}{\hat{\sigma}_R} \right) \left(\frac{\hat{\mathbf{b}}^\top}{\hat{\sigma}_R} \mathbf{D}_i \boldsymbol{\Omega} \frac{\hat{\mathbf{b}}}{\hat{\sigma}_R} \right) \end{aligned} \quad (30)$$

4 An ECME algorithm

The EM algorithm [3] is a general iterative algorithm for computing maximum likelihood estimates in the presence of missing data. For linear mixed-effects models we formulate an EM algorithm by considering the random effects \mathbf{b} to be unobserved data. In the terminology of the EM algorithm, we call the observed data, \mathbf{y} , the *incomplete* data, and \mathbf{y} augmented by \mathbf{b} , the *complete* data.

The EM algorithm has two steps: in the E step we compute Q , the expected log-likelihood (or deviance) for the complete data, and, in the M step, we maximize the expected log-likelihood (or minimize the expected deviance) with respect to the parameters in the model.

Liu and Rubin [4] derived the EM algorithm for linear mixed-effects models using \mathbf{b} as the missing data. In same paper they introduced *expectation conditional maximization either* (ECME) algorithms, which are an extension of the EM algorithm. In ECME algorithms the M step is broken down into a number of conditional maximization steps and in each conditional maximization step either the original log-likelihood, ℓ , or its conditional expectation, Q , is maximized. The maximization in each step is done by placing constraints on the parameters in such a way that the collection of all the maximization steps

is with respect to the full parameter space.

In describing an EM algorithm we must distinguish between current values of parameters and updated values. We denote the current values of the parameters by $\boldsymbol{\beta}_0$, σ_0^2 and $\boldsymbol{\theta}_0$. These are either starting values or values obtained from the last E and M steps. The parameter estimates to be obtained after an E and an M step are $\boldsymbol{\beta}_1$, σ_1^2 and $\boldsymbol{\theta}_1$. The log-likelihood for the complete data is

$$-2\ell(\boldsymbol{\beta}, \sigma^2, \boldsymbol{\theta} | \mathbf{y}, \mathbf{b}) = (n + q) \log(2\pi\sigma^2) - \log |\boldsymbol{\Omega}| + \frac{\|\mathbf{y} - \mathbf{X}\boldsymbol{\beta} - \mathbf{Z}\mathbf{b}\|^2 + \mathbf{b}^\top \boldsymbol{\Omega} \mathbf{b}}{\sigma^2} \quad (31)$$

Because we can easily calculate $\widehat{\sigma^2}(\boldsymbol{\theta})$ and $\widehat{\boldsymbol{\beta}}(\boldsymbol{\theta})$ we define an update step in an ECME algorithm to be:

- (1) Given $\boldsymbol{\theta}_0$, set $\boldsymbol{\beta}_1 = \widehat{\boldsymbol{\beta}}(\boldsymbol{\theta}_0)$ and $\sigma_1^2 = \widehat{\sigma^2}(\boldsymbol{\theta}_0)$. Then the conditional distribution of \mathbf{b} is $\mathcal{N}\left(\widehat{\mathbf{b}}(\boldsymbol{\theta}), \sigma_1^2 (\mathbf{Z}^\top \mathbf{Z} + \boldsymbol{\Omega})^{-1}\right)$.
- (2) Choose $\boldsymbol{\theta}_1$ to minimize the conditional expectation of -2ℓ

$$\begin{aligned} Q(\boldsymbol{\theta} | \mathbf{y}, \sigma_1^2, \boldsymbol{\beta}_1, \boldsymbol{\theta}_0) &= \mathbb{E}_{\mathbf{b} | \boldsymbol{\theta}_0} [-2\ell(\boldsymbol{\beta}_1, \sigma_1^2, \boldsymbol{\theta} | \mathbf{y}, \mathbf{b})] \\ &= \mathbb{E}_{\mathbf{b} | \boldsymbol{\theta}_0} \left[c - \log |\boldsymbol{\Omega}| + \mathbf{b}^\top \boldsymbol{\Omega} \mathbf{b} / \sigma_1^2 \right] \\ &= c - \log |\boldsymbol{\Omega}| + \frac{\widehat{\mathbf{b}}^\top \boldsymbol{\Omega} \widehat{\mathbf{b}}}{\sigma_1^2} + \text{tr} \left[\boldsymbol{\Omega} (\mathbf{Z}^\top \mathbf{Z} + \boldsymbol{\Omega}(\boldsymbol{\theta}_0))^{-1} \right] \end{aligned} \quad (32)$$

Thus $\boldsymbol{\theta}_1$ satisfies

$$\nabla_{\boldsymbol{\theta}} Q = \text{tr} \left[\mathbf{D} \boldsymbol{\Omega} \left(\frac{\widehat{\mathbf{b}}(\boldsymbol{\theta}_0) \widehat{\mathbf{b}}(\boldsymbol{\theta}_0)^\top}{\sigma_1} \frac{1}{\sigma_1} + (\mathbf{Z}^\top \mathbf{Z} + \boldsymbol{\Omega}(\boldsymbol{\theta}_0))^{-1} - \boldsymbol{\Omega}(\boldsymbol{\theta}_1)^{-1} \right) \right] = \mathbf{0} \quad (33)$$

A similar derivation shows that, for the REML criterion, $\boldsymbol{\theta}_1$ satisfies

$$\nabla_{\boldsymbol{\theta}} Q_R = \text{tr} \left[\mathbf{D} \boldsymbol{\Omega} \left(\frac{\widehat{\mathbf{b}}(\boldsymbol{\theta}_0) \widehat{\mathbf{b}}(\boldsymbol{\theta}_0)^\top}{\sigma_R} \frac{1}{\sigma_R} + \mathbf{V}_b(\boldsymbol{\theta}_0) - \boldsymbol{\Omega}(\boldsymbol{\theta}_1)^{-1} \right) \right] = \mathbf{0} \quad (34)$$

From the similarity of (33) to (27) (and of (34) to (28)) we can see that a stationary point of this ECME algorithm will be a critical value of the profiled log-likelihood.

5 Computational methods for a single grouping factor

The results in the previous sections provide concise expressions for the profiled log-likelihood (10), the profiled log-restricted-likelihood (15); the ECME

increments (33 or 34); and the gradient (27 or 28) and Hessian (29 or 30) of the profiled objective functions. All these expressions depend on being able to evaluate the initial decomposition (18) and, for several different $\boldsymbol{\theta}$, the decomposition (19), which can be a formidable computational problem because the matrices \mathbf{Z} and $\boldsymbol{\Omega}$ can be very large. However, these matrices generally are sparse and, by exploiting the sparsity, we can provide computationally feasible methods for all these results.

The sparsity in \mathbf{Z} and $\boldsymbol{\Omega}$ occurs when the random effects vector \mathbf{b} is divided into small components associated with one or more factors that group the observations. In the simplest situation there is one grouping factor, or one set of *experimental units*, and the model can be written

$$\mathbf{y}_i = \mathbf{X}_i\boldsymbol{\beta} + \mathbf{Z}_i\mathbf{b}_i + \boldsymbol{\epsilon}_i, \quad \mathbf{b}_i \sim \mathcal{N}(\mathbf{0}, \sigma^2\boldsymbol{\Omega}_1^{-1}), \quad \boldsymbol{\epsilon}_i \sim \mathcal{N}(\mathbf{0}, \sigma^2\mathbf{I}), \quad i = 1, \dots, m, \\ \boldsymbol{\epsilon}_i \perp \boldsymbol{\epsilon}_j, \quad \mathbf{b}_i \perp \mathbf{b}_j, \quad i \neq j; \quad \boldsymbol{\epsilon}_i \perp \mathbf{b}_j, \quad \text{all } i, j \quad (35)$$

where \mathbf{y}_i is the vector of length n_i of responses for unit i ; \mathbf{X}_i is the $n_i \times p$ fixed-effects model matrix for unit i ; and \mathbf{Z}_i is the $n_i \times q_1$ model matrix for unit i and the random effects \mathbf{b}_i corresponding to that unit.

Because we only have one grouping factor in this model we say that we have one “level” of random effects. However, this terminology is not universal. In particular, this model is called the “two-level” model in the multilevel modeling literature (e.g. [5]) because it has two levels of random variation.

To convert model (35) to the form (1) we would set

$$\mathbf{b} = \begin{bmatrix} \mathbf{b}_1 \\ \mathbf{b}_2 \\ \vdots \\ \mathbf{b}_{m_1} \end{bmatrix}, \quad \mathbf{y} = \begin{bmatrix} \mathbf{y}_1 \\ \mathbf{y}_2 \\ \vdots \\ \mathbf{y}_{m_1} \end{bmatrix}, \quad \mathbf{X} = \begin{bmatrix} \mathbf{X}_1 \\ \mathbf{X}_2 \\ \vdots \\ \mathbf{X}_{m_1} \end{bmatrix}, \quad \mathbf{Z} = \begin{bmatrix} \mathbf{Z}_1 & \mathbf{0} & \dots & \mathbf{0} \\ \mathbf{0} & \mathbf{Z}_2 & \dots & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{0} & \mathbf{0} & \dots & \mathbf{Z}_{m_1} \end{bmatrix}, \quad (36)$$

and

$$\boldsymbol{\Omega} = \begin{bmatrix} \boldsymbol{\Omega}_1 & \mathbf{0} & \dots & \mathbf{0} \\ \mathbf{0} & \boldsymbol{\Omega}_1 & \dots & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{0} & \mathbf{0} & \dots & \boldsymbol{\Omega}_1 \end{bmatrix} = \mathbf{I}_{m_1} \otimes \boldsymbol{\Omega}_1 \quad (37)$$

where \mathbf{I}_{m_1} is the $m_1 \times m_1$ identity matrix and \otimes is the Kronecker product. That is, $\boldsymbol{\Omega}$ is a block-diagonal matrix whose diagonal is m_1 copies of $\boldsymbol{\Omega}_1$. Similarly, a relative precision factor is $\boldsymbol{\Delta} = \mathbf{I}_{m_1} \otimes \boldsymbol{\Delta}_1$ where $\boldsymbol{\Delta}_1$ is any $q_1 \times q_1$ matrix satisfying $\boldsymbol{\Delta}_1^\top \boldsymbol{\Delta}_1 = \boldsymbol{\Omega}_1$.

The calculation of \mathbf{R}_{XX} ,

$$\mathbf{R}_{ZZ} = \begin{bmatrix} \mathbf{R}_{ZZ(1)} & \mathbf{0} & \dots & \mathbf{0} \\ \mathbf{0} & \mathbf{R}_{ZZ(2)} & \dots & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{0} & \mathbf{0} & \dots & \mathbf{R}_{ZZ(m_1)} \end{bmatrix}, \quad \mathbf{R}_{ZX} = \begin{bmatrix} \mathbf{R}_{ZX(1)} \\ \mathbf{R}_{ZX(2)} \\ \vdots \\ \mathbf{R}_{ZX(m_1)} \end{bmatrix} \quad (38)$$

and $\mathbf{r}_{Zy} = \left[\mathbf{r}_{Zy(1)}, \mathbf{r}_{Zy(2)}, \dots, \mathbf{r}_{Zy(m_1)} \right]^\top$ can be performed in blocks, using a series of QR decompositions based on the corresponding blocks in \mathbf{T}

$$\begin{bmatrix} \Delta_1 \\ \mathbf{T}_{ZZ(i)} \end{bmatrix} = \mathbf{Q}_i \begin{bmatrix} \mathbf{R}_{ZZ(i)} \\ \mathbf{0} \end{bmatrix} \text{ followed by } \mathbf{Q}_i^\top \begin{bmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{T}_{ZX(i)} & \mathbf{t}_{Zy(i)} \end{bmatrix} = \begin{bmatrix} \mathbf{R}_{ZX(i)} & \mathbf{r}_{Zy(i)} \\ \mathbf{R}_{XX(i)} & \mathbf{r}_{Xy(i)} \end{bmatrix}, \quad (39)$$

and

$$\begin{bmatrix} \mathbf{R}_{XX(1)} & \mathbf{r}_{Xy(1)} \\ \vdots & \vdots \\ \mathbf{R}_{XX(m_1)} & \mathbf{r}_{Xy(m_1)} \\ \mathbf{T}_{XX} & \mathbf{t}_{Xy} \\ \mathbf{0} & \mathbf{t}_{yy} \end{bmatrix} = \mathbf{Q}_0 \begin{bmatrix} \mathbf{R}_{XX} & \mathbf{r}_{Xy} \\ \mathbf{0} & r_{yy} \\ \mathbf{0} & \mathbf{0} \end{bmatrix}. \quad (40)$$

Noting that $|\Omega| = |\Omega_1|^{m_1} = |\Delta_1|^{2m_1}$ and that $|\mathbf{Z}^\top \mathbf{Z} + \Omega| = \prod_{i=1}^{m_1} |\mathbf{R}_{ZZ(i)}|^2$, we can evaluate the profiled log-likelihood (10) or the profiled log-restricted-likelihood (15).

To evaluate the ECME increments we note that

$$\begin{aligned} \text{tr} \left[\mathbf{D}_i(\Omega) \Omega^{-1} \right] &= m_1 \text{tr} \left[\mathbf{D}_i(\Omega_1) \Omega_1^{-1} \right] = m_1 \text{tr} \left[\Delta_1^{-\top} \mathbf{D}_i(\Omega_1) \Delta_1^{-1} \right] \\ \text{tr} \left[\mathbf{D}_i(\Omega) \left(\frac{\hat{\mathbf{b}}(\theta_0) \hat{\mathbf{b}}(\theta_0)^\top}{\sigma_1} + (\mathbf{Z}^\top \mathbf{Z} + \Omega(\theta_0))^{-1} \right) \right] &= \text{tr} \left[\mathbf{A}_1^\top \mathbf{D}_i(\Omega) \mathbf{A}_1 \right] \\ \text{tr} \left[\mathbf{D}_i(\Omega) \left(\frac{\hat{\mathbf{b}}(\theta_0) \hat{\mathbf{b}}(\theta_0)^\top}{\sigma_R} + \mathbf{V}_b \right) \right] &= \text{tr} \left[\mathbf{A}_{R1}^\top \mathbf{D}_i(\Omega) \mathbf{A}_{R1} \right] \end{aligned} \quad (41)$$

where the matrices \mathbf{A}_1 and \mathbf{A}_{R1} are obtained from the QR decompositions

$$\mathbf{U}_1 \mathbf{A}_1 = \begin{bmatrix} \hat{\mathbf{b}}_1^\top / \sigma_1 \\ \mathbf{R}_{ZZ(1)}^{-\top} \\ \vdots \\ \hat{\mathbf{b}}_{m_1}^\top / \sigma_1 \\ \mathbf{R}_{ZZ(m_1)}^{-\top} \end{bmatrix} \quad \text{and} \quad \mathbf{U}_{R1} \mathbf{A}_{R1} = \begin{bmatrix} \hat{\mathbf{b}}_1^\top / \sigma_R \\ \mathbf{R}_{ZZ(1)}^{-\top} \\ -\mathbf{R}_{XX}^{-\top} \mathbf{R}_{ZX(1)}^\top \mathbf{R}_{ZZ(1)}^{-\top} \\ \vdots \\ \hat{\mathbf{b}}_{m_1}^\top / \sigma_R \\ \mathbf{R}_{ZZ(m_1)}^{-\top} \\ -\mathbf{R}_{XX}^{-\top} \mathbf{R}_{ZX(m_1)}^\top \mathbf{R}_{ZZ(m_1)}^{-\top} \end{bmatrix} \quad (42)$$

The matrix $\mathbf{\Omega}_1$ must be positive definite and symmetric. If no further structure is imposed on it then $k = q_1(q_1 + 1)/2$ and a suitable value of $\mathbf{\Delta}_1(\boldsymbol{\theta}_1)$ can be calculated as $\mathbf{\Delta}_1 = \sqrt{m_1} \mathbf{A}_1^{-\top}$ for ML estimation and $\mathbf{\Delta}_{R1} = \sqrt{m_1} \mathbf{A}_{R1}^{-\top}$ for REML. If further structure is imposed on $\mathbf{\Omega}_1(\boldsymbol{\theta})$, so that $k < q_1(q_1 + 1)/2$, then the gradient equations, (33) or (34), must be solved for the ECME update.

Evaluation of the gradients (27 or 28) of the profiled objective functions can be simplified in the same way as the ECME increment is, as can the first term in the Hessian (29 or 30).

The other terms in the Hessian can be simplified in various ways. Because $\mathbf{D}_i \mathbf{\Omega} = \mathbf{I}_{m_1} \otimes \mathbf{D}_i \mathbf{\Omega}_1$ and $\mathbf{\Omega}^{-1} = \mathbf{I}_{m_1} \otimes \mathbf{\Omega}_1^{-1}$

$$\begin{aligned} \text{tr} \left(\mathbf{D}_j \mathbf{\Omega} \mathbf{\Omega}^{-1} \mathbf{D}_i \mathbf{\Omega} \mathbf{\Omega}^{-1} \right) &= m_1 \text{tr} \left(\mathbf{D}_j \mathbf{\Omega}_1 \mathbf{\Omega}_1^{-1} \mathbf{D}_i \mathbf{\Omega}_1 \mathbf{\Omega}_1^{-1} \right) \\ &= m_1 \text{tr} \left(\mathbf{\Delta}_1^{-\top} \mathbf{D}_j \mathbf{\Omega}_1 \mathbf{\Delta}_1^{-1} \mathbf{\Delta}_1^{-\top} \mathbf{D}_i \mathbf{\Omega}_1 \mathbf{\Delta}_1^{-1} \right) \end{aligned} \quad (43)$$

From the vectors $\mathbf{c}_{\ell i}$, $\mathbf{d}_{\ell i}$, and $\mathbf{g}_{\ell i}$, $\ell = 1, \dots, m_1$; $i = 1, \dots, k$ defined as

$$\begin{aligned} \mathbf{c}_{\ell i} &= \mathbf{D}_i \mathbf{\Omega}_1 \hat{\mathbf{b}}_\ell \\ \mathbf{d}_{\ell i} &= \mathbf{R}_{ZZ(\ell)}^{-\top} \mathbf{c}_{\ell i} \\ \mathbf{g}_{\ell i} &= -\mathbf{R}_{XX}^{-\top} \mathbf{R}_{ZX(\ell)}^\top \mathbf{d}_{\ell i} \end{aligned} \quad (44)$$

we can evaluate

$$\begin{aligned} \left(\hat{\mathbf{b}}^\top \mathbf{D}_j \mathbf{\Omega} \hat{\mathbf{b}} \right) \left(\hat{\mathbf{b}}^\top \mathbf{D}_i \mathbf{\Omega} \hat{\mathbf{b}} \right) &= \left(\sum_{\ell=1}^{m_1} \mathbf{b}_\ell^\top \mathbf{c}_{\ell j} \right) \left(\sum_{\ell=1}^{m_1} \mathbf{b}_\ell^\top \mathbf{c}_{\ell i} \right) \\ \hat{\mathbf{b}}^\top \mathbf{D}_j \mathbf{\Omega} \mathbf{V}_b \mathbf{D}_i \mathbf{\Omega} \hat{\mathbf{b}} &= \sum_{\ell=1}^{m_1} \mathbf{d}_{\ell j}^\top \mathbf{d}_{\ell i} + \left(\sum_{\ell=1}^{m_1} \mathbf{g}_{\ell j} \right)^\top \left(\sum_{\ell=1}^{m_1} \mathbf{g}_{\ell i} \right) \end{aligned} \quad (45)$$

The remaining terms in the Hessian expressions are evaluated as

$$\begin{aligned} \text{tr} \left[\mathbf{D}_j \boldsymbol{\Omega} (\mathbf{Z}^\top \mathbf{Z} + \boldsymbol{\Omega})^{-1} \mathbf{D}_i \boldsymbol{\Omega} (\mathbf{Z}^\top \mathbf{Z} + \boldsymbol{\Omega})^{-1} \right] = \\ \sum_{\ell=1}^{m_1} \text{tr} \left[\mathbf{R}_{ZZ(\ell)}^{-\top} \mathbf{D}_j \boldsymbol{\Omega}_1 \mathbf{R}_{ZZ(\ell)}^{-1} \mathbf{R}_{ZZ(\ell)}^{-\top} \mathbf{D}_i \boldsymbol{\Omega}_1 \mathbf{R}_{ZZ(\ell)}^{-1} \right] \end{aligned} \quad (46)$$

or

$$\begin{aligned} \text{tr} [\mathbf{D}_j \boldsymbol{\Omega} \mathbf{V}_b \mathbf{D}_i \boldsymbol{\Omega} \mathbf{V}_b] = \\ \sum_{\ell=1}^{m_1} \text{tr} \left[\mathbf{R}_{ZZ(\ell)}^{-\top} \mathbf{D}_j \boldsymbol{\Omega}_1 \mathbf{R}_{ZZ(\ell)}^{-1} \mathbf{R}_{ZZ(\ell)}^{-\top} \mathbf{D}_i \boldsymbol{\Omega}_1 \mathbf{R}_{ZZ(\ell)}^{-1} \right] \\ + 2 \sum_{\ell=1}^{m_1} \text{tr} \left[\mathbf{R}_{ZZ(\ell)}^{-\top} \mathbf{D}_j \boldsymbol{\Omega}_1 \mathbf{R}_{ZZ(\ell)}^{-1} \mathbf{R}_{ZX(\ell)} \mathbf{R}_{XX}^{-1} \mathbf{R}_{XX}^{-\top} \mathbf{R}_{ZX(\ell)}^\top \mathbf{R}_{ZZ(\ell)}^{-\top} \mathbf{D}_i \boldsymbol{\Omega}_1 \mathbf{R}_{ZZ(\ell)}^{-1} \right] \\ + \text{tr} \left[\left(\sum_{\ell=1}^{m_1} \mathbf{R}_{XX}^{-\top} \mathbf{R}_{ZX(\ell)}^\top \mathbf{R}_{ZZ(\ell)}^{-\top} \mathbf{D}_j \boldsymbol{\Omega}_1 \mathbf{R}_{ZZ(\ell)}^{-1} \mathbf{R}_{ZX(\ell)} \mathbf{R}_{XX}^{-1} \right) \right. \\ \left. \left(\sum_{\ell=1}^{m_1} \mathbf{R}_{XX}^{-\top} \mathbf{R}_{ZX(\ell)}^\top \mathbf{R}_{ZZ(\ell)}^{-\top} \mathbf{D}_i \boldsymbol{\Omega}_1 \mathbf{R}_{ZZ(\ell)}^{-1} \mathbf{R}_{ZX(\ell)} \mathbf{R}_{XX}^{-1} \right) \right] \end{aligned} \quad (47)$$

By working with the components of $\boldsymbol{\Omega}$ instead of the whole matrix we save a considerable amount of storage and computation because q_1 is typically very small (values of one or two are common) while m_1 can be very large. We have worked with cases where m_1 is in the millions. In most cases $p \geq q_1$ but it is unusual for p to exceed, say, one hundred.

After the pre-decomposition the storage required is approximately $m_1 q_1 (q_1 + p)$ locations and the number of floating point operations (FLOPs) per iteration for the function evaluation is on the order of $m_1 q_1^3$ for the decompositions and $m_1 p q_1^2$ for the multiplications in (39), followed by $m_1 q_1 p^2$ for (40). Evaluation of the determinants of $\mathbf{R}_{ZZ(i)}$, $i = 1, \dots, m_1$ and \mathbf{R}_{XX} is trivial because these matrices are triangular. The dominant term in the FLOP count is $m_1 q_1 (q_1 + p)^2$.

This is also the order of the computation for the ECME update, the gradient calculation and the Hessian calculation.

We can simplify some of the expressions for the Hessian by noting that, for arbitrary $q_1 \times q_1$ matrices \mathbf{U} and \mathbf{V} ,

$$\begin{aligned} \text{tr} [\mathbf{D}_i \boldsymbol{\Omega}_1 \mathbf{U} \mathbf{D}_j \boldsymbol{\Omega}_1 \mathbf{V}] &= \sum_{k_1} \sum_{k_2} \sum_{k_3} \sum_{k_4} d_{k_1 k_2}^{(i)} d_{k_3 k_4}^{(j)} \text{tr} [\mathbf{e}_{k_1} \mathbf{e}_{k_2}^\top \mathbf{U} \mathbf{e}_{k_3} \mathbf{e}_{k_4}^\top \mathbf{V}] \\ &= \sum_{k_1} \sum_{k_2} \sum_{k_3} \sum_{k_4} d_{k_1 k_2}^{(i)} d_{k_3 k_4}^{(j)} \left(\mathbf{e}_{k_2}^\top \mathbf{U} \mathbf{e}_{k_3} \mathbf{e}_{k_4}^\top \mathbf{V} \mathbf{e}_{k_1} \right) \\ &= \sum_{k_1} \sum_{k_2} \sum_{k_3} \sum_{k_4} d_{k_1 k_2}^{(i)} d_{k_3 k_4}^{(j)} u_{k_2 k_3} v_{k_4 k_1} \end{aligned}$$

where $d_{jk}^{(i)}$ is the (j, k) th element of $\mathbf{D}_i \boldsymbol{\Omega}_1$, \mathbf{e}_k is the k -th column of the identity matrix of size $q_1 \times q_1$, and u_{ij} and v_{ij} are the (i, j) th elements of \mathbf{U} and \mathbf{V} , respectively.

For the one-level model, all the terms involved in the Hessian can be rewritten using this principle and so the full Hessian can be written as

$$\sum_{k_1 k_2 k_3 k_4} d_{k_1 k_2}^{(i)} d_{k_3 k_4}^{(j)} \mathbf{H}_{k_1 k_2 k_3 k_4}$$

where \mathbf{H} is an array of size $q_1 \times q_1 \times q_1 \times q_1$.

6 Models with multiple grouping factors

We will consider models with multiple grouping factors where the random effects associated with each grouping factor are independent between groups and are i.i.d within groups. That is, the matrix $\boldsymbol{\Omega}$ consists of s blocks of the form $\mathbf{I}_{m_j} \otimes \boldsymbol{\Omega}_j$, $j = 1, \dots, s$. We will give details for $s = 2$, in which case

$$\boldsymbol{\Omega} = \begin{bmatrix} \mathbf{I}_{m_1} \otimes \boldsymbol{\Omega}_1 & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_{m_2} \otimes \boldsymbol{\Omega}_2 \end{bmatrix}. \quad (48)$$

Extensions to more than two grouping factors follow naturally.

The block-diagonal structure in (48) will also be present in the factor, $\boldsymbol{\Delta}$, and in all the derivatives with respect to components of $\boldsymbol{\theta}$. That is, $\boldsymbol{\Delta}$ consists of s blocks of the form $\mathbf{I}_{m_j} \otimes \boldsymbol{\Delta}_j$, $j = 1, \dots, s$ where $\boldsymbol{\Delta}_j$ is of size $q_j \times q_j$. We assume that each component of $\boldsymbol{\theta}$ determines only one block in $\boldsymbol{\Omega}$ and we designate the block associated with component j as $b(j)$ where $1 \leq b(j) \leq s$, $j = 1, \dots, k$. Then $\mathbf{D}_j \boldsymbol{\Omega}$ is zero except for the $b(j)$ th diagonal block which is of the form $\mathbf{I}_{m_{b(j)}} \otimes \mathbf{D}_j \boldsymbol{\Omega}_{b(j)}$. Notice that the second derivative, $\mathbf{D}_j \mathbf{D}_i \boldsymbol{\Omega}$, will be zero if $b(i) \neq b(j)$.

We divide the random effects vector \mathbf{b} into $s = 2$ blocks and subdivide the j th block into m_j components of length q_j denoted by \mathbf{b}_{ji} , $j = 1, \dots, s$; $i = 1, \dots, m_j$.

If $\mathbf{Z}^\top \mathbf{Z}$ is split into blocks corresponding to the blocks in $\boldsymbol{\Omega}$ then the diagonal blocks in $\mathbf{Z}^\top \mathbf{Z}$ are themselves block diagonal but, unlike the situation with a single grouping factor, $\mathbf{Z}^\top \mathbf{Z}$ and \mathbf{T}_{ZZ} will have non-zero off-diagonal blocks. These off-diagonal blocks can be sparse or dense according to whether the grouping factors are nested or crossed. We distinguish three cases: completely crossed, partially crossed, and strictly nested.

Completely crossed grouping factors usually occur in designed experiments. For example, biological assays are often conducted by measuring the optical density of liquid samples in wells arranged in a grid on a plate. A common arrangement is 96 wells in a grid of 8 rows by 12 columns. If we assigned one set of random effects to the rows and another set of random effects to the columns then each row would occur with each column, resulting in completely crossed random effects.

An example of nested grouping factors would be a longitudinal study, say records of annual achievement test scores, of students in several schools. For the student factor to be strictly nested within the school factor we require that each student attend only one school during the period of the study. In most large studies this will not be the case. We expect some students will attend more than one school but we do not expect every student to attend every school. In such a case the grouping factors are neither strictly nested nor completely crossed. We describe this situation as partially crossed grouping factors.

If we divide $\mathbf{Z} = \begin{bmatrix} \mathbf{Z}_1 & \mathbf{Z}_2 \end{bmatrix}$ according to the two levels of random effects and correspondingly divide

$$\mathbf{T}_{ZZ} = \begin{bmatrix} \mathbf{T}_{ZZ11} & \mathbf{T}_{ZZ12} \\ \mathbf{0} & \mathbf{T}_{ZZ22} \end{bmatrix} \text{ and } \mathbf{R}_{ZZ} = \begin{bmatrix} \mathbf{R}_{ZZ11} & \mathbf{R}_{ZZ12} \\ \mathbf{0} & \mathbf{R}_{ZZ22} \end{bmatrix} \quad (49)$$

then the block diagonal structure of $\mathbf{Z}_1^\top \mathbf{Z}_1$ is also present in \mathbf{T}_{ZZ11} and we can store and manipulate it accordingly. Even though $\mathbf{Z}_2^\top \mathbf{Z}_2$ is block diagonal there can be non-zero off-diagonals induced in \mathbf{T}_{ZZ22} from \mathbf{T}_{ZZ12} when taking the Cholesky decomposition. This is called “fill-in” [6].

If we consider the matrix \mathbf{T}_{ZZ12} to consist of $m_1 \times m_2$ blocks of size $q_1 \times q_2$ then the (i,j)th such block will necessarily be zero if the i th level of the first grouping factor does not occur in combination with the j th level of the second grouping factor.

For strictly nested grouping factors there will be only one non-zero block in each of the m_1 sets of q_1 adjacent rows. Furthermore, this pattern guarantees that there will be no fill-in of \mathbf{T}_{ZZ22} and that it can be stored as a block-diagonal matrix consisting of m_2 blocks of size $q_2 \times q_2$. The matrix \mathbf{R}_{ZZ} has the same pattern of sparsity so \mathbf{T}_{ZZ} and \mathbf{R}_{ZZ} only require storage of size $(m_1 q_1 + m_2 q_2)(q_1 + q_2 + p)$ and computing \mathbf{R}_{ZZ} is of order $(m_1 q_1 + m_2 q_2)(q_1 + q_2 + p)^2$ FLOPs. The same sparsity pattern will be present in \mathbf{R}_{ZZ}^{-1} , which we can calculate explicitly, allowing for the evaluation of all the formulas for the ECME update, the gradient, and the Hessian. These results generalize to an arbitrary number of strictly nested levels of grouping factors (and are

incorporated this way in our software).

For fully crossed grouping factors \mathbf{T}_{ZZ12} and \mathbf{T}_{ZZ22} are both dense and would need to be stored and manipulated as such. Because we can condense the calculation for the first group of random effects only, we choose the order of the grouping factors so that $m_1q_1 \geq m_2q_2$. Once the order of the grouping factors is established then the calculations involving \mathbf{T}_{ZZ12} and \mathbf{T}_{ZZ22} are essentially the same as those involving \mathbf{T}_{ZX} and \mathbf{T}_{XX} except that producing \mathbf{R}_{ZZ22} from \mathbf{T}_{ZZ22} also involves $\mathbf{\Delta}_2$.

The amount of storage required for fully crossed grouping factors is on the order of $(m_1q_1 + m_2q_2)(q_1 + m_2q_2 + p)$ locations and producing \mathbf{R}_{ZZ} requires on the order of $(m_1q_1 + m_2q_2)(q_1 + m_2q_2 + p)^2$ FLOPs, which, obviously, can be considerably more than $(m_1q_1 + m_2q_2)(q_1 + q_2 + p)^2$. However, fully crossed grouping factors result in $n \geq m_1m_2$ and usually come from designed experiments so we do not expect m_2q_2 to be extremely large.

The most interesting case is partially crossed grouping factors where \mathbf{T}_{ZZ12} is sparse but does not obey that pattern that there is only one non-zero $q_1 \times q_2$ block in each of the m_1 sets of q_1 adjacent rows. Having more than one non-zero block in such a set of rows (say because a student attended more than one school during the course of the study) does not generate fill-in in \mathbf{T}_{ZZ12} but does generate fill-in in \mathbf{T}_{ZZ22} and in \mathbf{R}_{ZZ22} .

If the extent of the crossing is moderate (i.e. an individual student may attend more than one school during the study but no student attends a large proportion of all the schools in the study, so that most of the $q_1 \times q_2$ blocks in any one of the m_1 sets of q_1 adjacent rows are zero), then it will be advantageous to use sparse matrix representations of \mathbf{T}_{ZZ12} (and \mathbf{R}_{ZZ12}) and to generate \mathbf{R}_{ZZ22} using methods for the Cholesky decomposition of sparse semi-definite matrices. This particular calculation has been extensively studied because it is important in the implementation of interior-point methods in mathematical programming (Wright (1997, pp. 253–254) [6]).

A critical part of algorithms for the Cholesky decomposition of sparse semi-definite matrices is choosing an ordering of the columns of \mathbf{Z}_2 so as to minimize fill-in in \mathbf{R}_{ZZ22} . This only needs to be done once and can be based on the pattern of crossing of the grouping factors.

The amount of storage and computation required to generate \mathbf{R}_{ZZ} for partially crossed grouping factors will fall between that for strictly nested grouping factors and that for fully crossed grouping factors. It is unlikely that \mathbf{R}_{ZZ} can be inverted in place when the grouping factors are partially crossed. We expect that it will be feasible to evaluate the objective function, the ECME update, and the gradient but that it may not be feasible to evaluate the Hessian.

7 Implementation

We have implemented the computational methods for evaluating the profiled objective function (either log-likelihood or log-restricted-likelihood), the ECME increment and the gradient in the `lme4` package for R [7] (www.r-project.org). Some comparisons of the speed and stability of this implementation versus our previous implementation in the `nlme` package for R are available in Debroy (2003) [8].

Future versions of this software will allow for crossed random effects and will provide the Hessian, at least for the completely crossed and strictly nested cases.

8 Acknowledgements

This work was supported by U.S. Army Medical Research and Materiel Command under Contract No. DAMD17-02-C-0119. The views, opinions and/or findings contained in this report are those of the authors and should not be construed as an official Department of the Army position, policy or decision unless so designated by other documentation. We thank Deepayan Sarkar for helpful discussions.

A Differentiating the penalized RSS

Using (21) and the relationships

$$D\Phi = \begin{bmatrix} \mathbf{0} \\ \mathbf{I} \end{bmatrix} D\Delta \begin{bmatrix} \mathbf{I} & \mathbf{0} \end{bmatrix}, \quad \Phi^+ \tilde{\mathbf{y}} = \begin{bmatrix} \hat{\mathbf{b}} \\ \hat{\boldsymbol{\beta}} \end{bmatrix}, \quad \text{and} \quad P^\perp \tilde{\mathbf{y}} = \begin{bmatrix} \mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}} - \mathbf{Z}\hat{\mathbf{b}} \\ -\Delta\hat{\mathbf{b}} \end{bmatrix} \quad (\text{A.1})$$

we obtain

$$\nabla r_{yy}^2(\boldsymbol{\theta}) = 2\hat{\mathbf{b}}^\top \Delta^\top D\Delta\hat{\mathbf{b}} = \hat{\mathbf{b}}^\top D\Omega\hat{\mathbf{b}} \quad (\text{A.2})$$

The (j, i) th element of the Hessian $\nabla^2 r_{yy}^2(\boldsymbol{\theta}) = D_j D_i r_{yy}^2$ is

$$\begin{aligned}
\frac{1}{2} D_j D_i r_{yy}^2 &= -D_j \left(\tilde{\mathbf{y}}^\top P^\perp D_i \Phi \Phi^+ \tilde{\mathbf{y}} \right) \\
&= -\tilde{\mathbf{y}}^\top \left(D_j P^\perp D_i \Phi \Phi^+ + P^\perp D_j D_i \Phi \Phi^+ + P^\perp D_i \Phi D_j \Phi^+ \right) \tilde{\mathbf{y}} \\
&= \tilde{\mathbf{y}}^\top \left[P^\perp D_j \Phi \Phi^+ + (\Phi^+)^T D_j \Phi^\top P^\perp \right] D_i \Phi \Phi^+ \tilde{\mathbf{y}} + \hat{\mathbf{b}}^\top \Delta^\top D_j D_i \Delta \hat{\mathbf{b}} \\
&\quad + \tilde{\mathbf{y}}^\top P^\perp D_i \Phi \left[\Phi^+ D_j \Phi \Phi^+ - (\Phi^\top \Phi)^{-1} D_j \Phi^\top P^\perp \right] \tilde{\mathbf{y}} \\
&= \tilde{\mathbf{y}}^\top (\Phi^+)^T D_j \Phi^\top D_i \Phi \Phi^+ \tilde{\mathbf{y}} + \hat{\mathbf{b}}^\top \Delta^\top D_j D_i \Delta \hat{\mathbf{b}} \\
&\quad - \tilde{\mathbf{y}}^\top (\Phi^+)^T D_j \Phi^\top \Phi (\Phi^\top \Phi)^{-1} \Phi^\top D_i \Phi \Phi^+ \tilde{\mathbf{y}} \\
&\quad + \tilde{\mathbf{y}} P^\perp D_j \Phi (\Phi^\top \Phi)^{-1} \Phi^\top D_i \Phi \Phi^+ \tilde{\mathbf{y}} \\
&\quad + \tilde{\mathbf{y}} P^\perp D_i \Phi (\Phi^\top \Phi)^{-1} \Phi^\top D_j \Phi \Phi^+ \tilde{\mathbf{y}} \\
&\quad - \tilde{\mathbf{y}} P^\perp D_i \Phi (\Phi^\top \Phi)^{-1} D_j \Phi^\top P^\perp \tilde{\mathbf{y}}
\end{aligned} \tag{A.3}$$

Writing

$$\mathbf{w}_i = D_i \Phi \Phi^+ \tilde{\mathbf{y}} = \begin{bmatrix} \mathbf{0} \\ D_i \Delta \hat{\mathbf{b}} \end{bmatrix}, \tag{A.4}$$

$$\mathbf{u}_i = (\mathbf{R}^{-1})^\top \Phi^\top D_i \Phi \Phi^+ \tilde{\mathbf{y}} = (\mathbf{R}^{-1})^\top \begin{bmatrix} \Delta^\top D_i \Delta \hat{\mathbf{b}} \\ \mathbf{0} \end{bmatrix}, \tag{A.5}$$

and

$$\mathbf{v}_i = -(\mathbf{R}^{-1})^\top D_i \Phi^\top P^\perp \tilde{\mathbf{y}} = (\mathbf{R}^{-1})^\top \begin{bmatrix} D_i \Delta^\top \Delta \hat{\mathbf{b}} \\ \mathbf{0} \end{bmatrix}, \tag{A.6}$$

where \mathbf{R} is the Cholesky decomposition of $\Phi^\top \Phi = \mathbf{R}^\top \mathbf{R}$, expression (A.3) becomes

$$\frac{1}{2} D_j D_i r_{yy}^2(\boldsymbol{\theta}) = \hat{\mathbf{b}}^\top \Delta^\top D_j D_i \Delta \hat{\mathbf{b}} + \mathbf{w}_j^\top \mathbf{w}_i - (\mathbf{u}_j + \mathbf{v}_j)^\top (\mathbf{u}_i + \mathbf{v}_i) \tag{A.7}$$

Noting that $\mathbf{u}_i + \mathbf{v}_i = (\mathbf{R}^{-1})^\top \begin{bmatrix} D_i \Omega \hat{\mathbf{b}} \\ \mathbf{0} \end{bmatrix}$, we have

$$D_j D_i r_{yy}^2(\boldsymbol{\theta}) = \hat{\mathbf{b}}^\top [D_j D_i \Omega - D_j \Omega V_b D_i \Omega - D_i \Omega V_b D_j \Omega] \hat{\mathbf{b}}. \tag{A.8}$$

References

- [1] C. Henderson, *Applications of Linear Models in Animal Breeding*, University of Guelph, 1984.
- [2] G. H. Golub, V. Pereyra, The differentiation of pseudo-inverses and nonlinear least squares problems whose variables separate, *SIAM J. of Numerical Analysis* 10 (2) (1973) 413–432.
- [3] A. P. Dempster, N. M. Laird, D. B. Rubin, Maximum likelihood from incomplete data via the EM algorithm, *J. of the Royal Statistical Society, Series B* 39 (1977) 1–22.
- [4] C. Liu, D. B. Rubin, The ECME algorithm: A simple extension of EM and ECM with faster monotone convergence, *Biometrika* 81 (1994) 633–648.
- [5] S. W. Raudenbush, A. S. Bryk, *Hierarchical Linear Models: Applications and Data Analysis Methods*, 2nd Edition, Sage, 2002.
- [6] S. J. Wright, *Primal-Dual Interior Point Methods*, SIAM, 1997.
- [7] R. Ihaka, R. Gentleman, R: A language for data analysis and graphics, *Journal of Computational and Graphical Statistics* 5 (3) (1996) 299–314.
- [8] S. DebRoy, *Computational methods for mixed-effects models*, Ph.D. thesis, U. of Wisconsin – Madison (2003).