

# Computational methods for mixed models

Douglas Bates  
Department of Statistics  
University of Wisconsin – Madison

March 14, 2014

## Abstract

The `lme4` package provides R functions to fit and analyze several different types of mixed-effects models, including linear mixed models, generalized linear mixed models and nonlinear mixed models. In this vignette we describe the formulation of these models and the computational approach used to evaluate or approximate the log-likelihood of a model/data/parameter value combination.

## 1 Introduction

The `lme4` package provides R functions to fit and analyze linear mixed models, generalized linear mixed models and nonlinear mixed models. These models are called *mixed-effects models* or, more simply, *mixed models* because they incorporate both *fixed-effects* parameters, which apply to an entire population or to certain well-defined and repeatable subsets of a population, and *random effects*, which apply to the particular experimental units or observational units in the study. Such models are also called *multilevel* models because the random effects represent levels of variation in addition to the per-observation noise term that is incorporated in common statistical models such as linear regression models, generalized linear models and nonlinear regression models.

We begin by describing common properties of these mixed models and the general computational approach used in the `lme4` package. The estimates of the parameters in a mixed model are determined as the values that optimize

an objective function — either the likelihood of the parameters given the observed data, for maximum likelihood (ML) estimates, or a related objective function called the REML criterion. Because this objective function must be evaluated at many different values of the model parameters during the optimization process, we focus on the evaluation of the objective function and a critical computation in this evaluation — determining the solution to a penalized, weighted least squares (PWLS) problem.

The dimension of the solution of the PWLS problem can be very large, perhaps in the millions. Furthermore, such problems must be solved repeatedly during the optimization process to determine parameter estimates. The whole approach would be infeasible were it not for the fact that the matrices determining the PWLS problem are sparse and we can use sparse matrix storage formats and sparse matrix computations (Davis, 2006). In particular, the whole computational approach hinges on the extraordinarily efficient methods for determining the Cholesky decomposition of sparse, symmetric, positive-definite matrices embodied in the CHOLMOD library of C functions (Davis, 2005).

In the next section we describe the general form of the mixed models that can be represented in the `lme4` package and the computational approach embodied in the package. In the following section we describe a particular form of mixed model, called a linear mixed model, and the computational details for those models. In the fourth section we describe computational methods for generalized linear mixed models, nonlinear mixed models and generalized nonlinear mixed models.

## 2 Formulation of mixed models

A mixed-effects model incorporates two vector-valued random variables: the  $n$ -dimensional response vector,  $\mathbf{y}$ , and the  $q$ -dimensional random effects vector,  $\mathbf{b}$ . We observe the value,  $y$ , of  $\mathbf{y}$ . We do not observe the value of  $\mathbf{b}$ .

The random variable  $\mathbf{y}$  may be continuous or discrete. That is, the observed data,  $y$ , may be on a continuous scale or they may be on a discrete scale, such as binary responses or responses representing a count. In our formulation, the random variable  $\mathbf{b}$  is always continuous.

We specify a mixed model by describing the unconditional distribution of  $\mathbf{b}$  and the conditional distribution ( $\mathbf{y}|\mathbf{b} = \mathbf{b}$ ).

## 2.1 The unconditional distribution of $\mathcal{B}$

In our formulation, the unconditional distribution of  $\mathcal{B}$  is always a  $q$ -dimensional multivariate Gaussian (or “normal”) distribution with mean  $\mathbf{0}$  and with a parameterized covariance matrix,

$$\mathcal{B} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \Lambda(\boldsymbol{\theta}) \Lambda'(\boldsymbol{\theta})). \quad (1)$$

The scalar,  $\sigma$ , in (1), is called the *common scale parameter*. As we will see later, not all types of mixed models incorporate this parameter. We will include  $\sigma^2$  in the general form of the unconditional distribution of  $\mathcal{B}$  with the understanding that, in some models,  $\sigma \equiv 1$ .

The  $q \times q$  matrix  $\Lambda(\boldsymbol{\theta})$ , which is a left factor of the covariance matrix (when  $\sigma = 1$ ) or the relative covariance matrix (when  $\sigma \neq 1$ ), depends on an  $m$ -dimensional parameter  $\boldsymbol{\theta}$ . Typically  $m \ll q$ ; in the examples we show below it is always the case that  $m < 5$ , even when  $q$  is in the thousands. The fact that  $m$  is very small is important because, as we shall see, determining the parameter estimates in a mixed model can be expressed as an optimization problem with respect to  $\boldsymbol{\theta}$  only.

The parameter  $\boldsymbol{\theta}$  may be, and typically is, subject to constraints. For ease of computation, we require that the constraints be expressed as “box” constraints of the form  $\theta_{iL} \leq \theta_i \leq \theta_{iU}, i = 1, \dots, m$  for constants  $\theta_{iL}$  and  $\theta_{iU}, i = 1, \dots, m$ . We shall write the set of such constraints as  $\boldsymbol{\theta}_L \leq \boldsymbol{\theta} \leq \boldsymbol{\theta}_R$ . The matrix  $\Lambda(\boldsymbol{\theta})$  is required to be non-singular (i.e. invertible) when  $\boldsymbol{\theta}$  is not on the boundary.

## 2.2 The conditional distribution, $(\mathcal{Y}|\mathcal{B} = \mathbf{b})$

The conditional distribution,  $(\mathcal{Y}|\mathcal{B} = \mathbf{b})$ , must satisfy:

1. The conditional mean,  $\boldsymbol{\mu}_{\mathcal{Y}|\mathcal{B}}(\mathbf{b}) = \mathbb{E}[\mathcal{Y}|\mathcal{B} = \mathbf{b}]$ , depends on  $\mathbf{b}$  only through the value of the *linear predictor*,  $\mathbf{Z}\mathbf{b} + \mathbf{X}\boldsymbol{\beta}$ , where  $\boldsymbol{\beta}$  is the  $p$ -dimensional *fixed-effects* parameter vector and the *model matrices*,  $\mathbf{Z}$  and  $\mathbf{X}$ , are fixed matrices of the appropriate dimension. That is, the two model matrices must have the same number of rows and must have  $q$  and  $p$  columns, respectively. The number of rows in  $\mathbf{Z}$  and  $\mathbf{X}$  is a multiple of  $n$ , the dimension of  $\mathbf{y}$ .
2. The scalar distributions,  $(\mathcal{Y}_i|\mathcal{B} = \mathbf{b}), i = 1, \dots, n$ , all have the same form and are completely determined by the conditional mean,  $\boldsymbol{\mu}_{\mathcal{Y}|\mathcal{B}}(\mathbf{b})$

and, at most, one additional parameter,  $\sigma$ , which is the common scale parameter.

3. The scalar distributions,  $(\mathcal{Y}_i|\mathbf{B} = \mathbf{b}), i = 1, \dots, n$ , are independent. That is, the components of  $\mathbf{Y}$  are *conditionally independent* given  $\mathbf{B}$ .

An important special case of the conditional distribution is the multivariate Gaussian distribution of the form

$$(\mathbf{Y}|\mathbf{B} = \mathbf{b}) \sim \mathcal{N}(\mathbf{Z}\mathbf{b} + \mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{I}_n) \quad (2)$$

where  $\mathbf{I}_n$  denotes the identity matrix of size  $n$ . In this case the conditional mean,  $\boldsymbol{\mu}_{\mathbf{Y}|\mathbf{B}}(\mathbf{b})$ , is exactly the linear predictor,  $\mathbf{Z}\mathbf{b} + \mathbf{X}\boldsymbol{\beta}$ , a situation we will later describe as being an “identity link” between the conditional mean and the linear predictor. Models with conditional distribution (2) are called *linear mixed models*.

### 2.3 A change of variable to “spherical” random effects

Because the conditional distribution  $(\mathbf{Y}|\mathbf{B} = \mathbf{b})$  depends on  $\mathbf{b}$  only through the linear predictor, it is easy to express the model in terms of a linear transformation of  $\mathbf{B}$ . We define the linear transformation from a  $q$ -dimensional “spherical” Gaussian random variable,  $\mathbf{U}$ , to  $\mathbf{B}$  as

$$\mathbf{B} = \Lambda(\boldsymbol{\theta})\mathbf{U}, \quad \mathbf{U} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_q). \quad (3)$$

(The term “spherical” refers to the fact that contours of constant probability density for  $\mathbf{U}$  are spheres centered at the mean — in this case,  $\mathbf{0}$ .)

When  $\boldsymbol{\theta}$  is not on the boundary this is an invertible transformation. When  $\boldsymbol{\theta}$  is on the boundary the transformation can fail to be invertible. However, we will only need to be able to express  $\mathbf{B}$  in terms of  $\mathbf{U}$  and that transformation is well-defined, even when  $\boldsymbol{\theta}$  is on the boundary.

The linear predictor, as a function of  $\mathbf{u}$ , is

$$\boldsymbol{\gamma}(\mathbf{u}) = \mathbf{Z}\Lambda(\boldsymbol{\theta})\mathbf{u} + \mathbf{X}\boldsymbol{\beta}. \quad (4)$$

When we wish to emphasize the role of the model parameters,  $\boldsymbol{\theta}$  and  $\boldsymbol{\beta}$ , in the formulation of  $\boldsymbol{\gamma}$ , we will write the linear predictor as  $\boldsymbol{\gamma}(\mathbf{u}, \boldsymbol{\theta}, \boldsymbol{\beta})$ .

## 2.4 The conditional density ( $\mathcal{U}|\mathcal{Y} = \mathbf{y}$ )

Because we observe  $\mathbf{y}$  and do not observe  $\mathbf{b}$  or  $\mathbf{u}$ , the conditional distribution of interest, for the purposes of statistical inference, is ( $\mathcal{U}|\mathcal{Y} = \mathbf{y}$ ) (or, equivalently, ( $\mathcal{B}|\mathcal{Y} = \mathbf{y}$ )). This conditional distribution is always a continuous distribution with conditional probability density  $f_{\mathcal{U}|\mathcal{Y}}(\mathbf{u}|\mathbf{y})$ .

We can evaluate  $f_{\mathcal{U}|\mathcal{Y}}(\mathbf{u}|\mathbf{y})$ , up to a constant, as the product of the unconditional density,  $f_{\mathcal{U}}(\mathbf{u})$ , and the conditional density (or the probability mass function, whichever is appropriate),  $f_{\mathcal{Y}|\mathcal{U}}(\mathbf{y}|\mathbf{u})$ . We write this unnormalized conditional density as

$$h(\mathbf{u}|\mathbf{y}, \boldsymbol{\theta}, \boldsymbol{\beta}, \sigma) = f_{\mathcal{Y}|\mathcal{U}}(\mathbf{y}|\mathbf{u}, \boldsymbol{\theta}, \boldsymbol{\beta}, \sigma) f_{\mathcal{U}}(\mathbf{u}|\sigma). \quad (5)$$

We say that  $h$  is the “unnormalized” conditional density because all we know is that the conditional density is proportional to  $h(\mathbf{u}|\mathbf{y}, \boldsymbol{\theta}, \boldsymbol{\beta}, \sigma)$ . To obtain the conditional density we must normalize  $h$  by dividing by the value of the integral

$$L(\boldsymbol{\theta}, \boldsymbol{\beta}, \sigma|\mathbf{y}) = \int_{\mathbb{R}^q} h(\mathbf{u}|\mathbf{y}, \boldsymbol{\theta}, \boldsymbol{\beta}, \sigma) d\mathbf{u}. \quad (6)$$

We write the value of the integral (6) as  $L(\boldsymbol{\theta}, \boldsymbol{\beta}, \sigma|\mathbf{y})$  because it is exactly the *likelihood* of the parameters  $\boldsymbol{\theta}$ ,  $\boldsymbol{\beta}$  and  $\sigma$ , given the observed data  $\mathbf{y}$ . The *maximum likelihood (ML) estimates* of these parameters are the values that maximize  $L$ .

## 2.5 Determining the ML estimates

The general problem of maximizing  $L(\boldsymbol{\theta}, \boldsymbol{\beta}, \sigma|\mathbf{y})$  with respect to  $\boldsymbol{\theta}$ ,  $\boldsymbol{\beta}$  and  $\sigma$  can be formidable because each evaluation of this function involves a potentially high-dimensional integral and because the dimension of  $\boldsymbol{\beta}$  can be large. However, this general optimization problem can be split into manageable subproblems. Given a value of  $\boldsymbol{\theta}$  we can determine the *conditional mode*,  $\tilde{\mathbf{u}}(\boldsymbol{\theta})$ , of  $\mathbf{u}$  and the *conditional estimate*,  $\tilde{\boldsymbol{\beta}}(\boldsymbol{\theta})$  simultaneously using *penalized, iteratively re-weighted least squares* (PIRLS). The conditional mode and the conditional estimate are defined as

$$\begin{bmatrix} \tilde{\mathbf{u}}(\boldsymbol{\theta}) \\ \tilde{\boldsymbol{\beta}}(\boldsymbol{\theta}) \end{bmatrix} = \arg \max_{\mathbf{u}, \boldsymbol{\beta}} h(\mathbf{u}|\mathbf{y}, \boldsymbol{\theta}, \boldsymbol{\beta}, \sigma). \quad (7)$$

(It may look as if we have missed the dependence on  $\sigma$  on the left-hand side but it turns out that the scale parameter does not affect the location of the optimal values of quantities in the linear predictor.)

As is common in such optimization problems, we re-express the conditional density on the *deviance scale*, which is negative twice the logarithm of the density, where the optimization becomes

$$\begin{bmatrix} \tilde{\mathbf{u}}(\boldsymbol{\theta}) \\ \tilde{\boldsymbol{\beta}}(\boldsymbol{\theta}) \end{bmatrix} = \arg \min_{\mathbf{u}, \boldsymbol{\beta}} -2 \log (h(\mathbf{u}|\mathbf{y}, \boldsymbol{\theta}, \boldsymbol{\beta}, \sigma)). \quad (8)$$

It is this optimization problem that can be solved quite efficiently using PIRLS. In fact, for linear mixed models, which are described in the next section,  $\tilde{\mathbf{u}}(\boldsymbol{\theta})$  and  $\tilde{\boldsymbol{\beta}}(\boldsymbol{\theta})$  can be directly evaluated.

The second-order Taylor series expansion of  $-2 \log h$  at  $\tilde{\mathbf{u}}(\boldsymbol{\theta})$  and  $\tilde{\boldsymbol{\beta}}(\boldsymbol{\theta})$  provides the Laplace approximation to the profiled deviance. Optimizing this function with respect to  $\boldsymbol{\theta}$  provides the ML estimates of  $\boldsymbol{\theta}$ , from which the ML estimates of  $\boldsymbol{\beta}$  and  $\sigma$  (if used) are derived.

### 3 Methods for linear mixed models

As indicated in the introduction, a critical step in our methods for determining the maximum likelihood estimates of the parameters in a mixed model is solving a penalized, weighted least squares (PWLS) problem. We will motivate the general form of the PWLS problem by first considering computational methods for linear mixed models that result in a penalized least squares (PLS) problem.

Recall from §2.2 that, in a linear mixed model, both the conditional distribution,  $(\mathbf{y}|\mathbf{u} = \mathbf{u})$ , and the unconditional distribution,  $\mathbf{u}$ , are spherical Gaussian distributions and that the conditional mean,  $\boldsymbol{\mu}_{\mathbf{y}|\mathbf{u}}(\mathbf{u})$ , is the linear predictor,  $\boldsymbol{\gamma}(\mathbf{u})$ . Because all the distributions determining the model are continuous distributions, we consider their densities. On the deviance scale these are

$$\begin{aligned} -2 \log(f_{\mathbf{u}}(\mathbf{u})) &= q \log(2\pi\sigma^2) + \frac{\|\mathbf{u}\|^2}{\sigma^2} \\ -2 \log(f_{\mathbf{y}|\mathbf{u}}(\mathbf{y}|\mathbf{u})) &= n \log(2\pi\sigma^2) + \frac{\|\mathbf{y} - \mathbf{Z}\boldsymbol{\Lambda}(\boldsymbol{\theta})\mathbf{u} - \mathbf{X}\boldsymbol{\beta}\|^2}{\sigma^2} \\ -2 \log(h(\mathbf{u}|\mathbf{y}, \boldsymbol{\theta}, \boldsymbol{\beta}, \sigma)) &= (n + q) \log(2\pi\sigma^2) + \frac{\|\mathbf{y} - \boldsymbol{\gamma}(\mathbf{u}, \boldsymbol{\theta}, \boldsymbol{\beta})\|^2 + \|\mathbf{u}\|^2}{\sigma^2} \\ &= (n + q) \log(2\pi\sigma^2) + \frac{d(\mathbf{u}|\mathbf{y}, \boldsymbol{\theta}, \boldsymbol{\beta})}{\sigma^2} \end{aligned} \quad (9)$$

In (9) the *discrepancy* function,

$$d(\mathbf{u}|\mathbf{y}, \boldsymbol{\theta}, \boldsymbol{\beta}) = \|\mathbf{y} - \boldsymbol{\gamma}(\mathbf{u}, \boldsymbol{\theta}, \boldsymbol{\beta})\|^2 + \|\mathbf{u}\|^2 \quad (10)$$

has the form of a penalized residual sum of squares in that the first term,  $\|\mathbf{y} - \boldsymbol{\gamma}(\mathbf{u}, \boldsymbol{\theta}, \boldsymbol{\beta})\|^2$  is the residual sum of squares for  $\mathbf{y}$ ,  $\mathbf{u}$ ,  $\boldsymbol{\theta}$  and  $\boldsymbol{\beta}$  and the second term,  $\|\mathbf{u}\|^2$ , is a penalty on the size of  $\mathbf{u}$ . Notice that the discrepancy does not depend on the common scale parameter,  $\sigma$ .

### 3.1 The canonical form of the discrepancy

Using a so-called “pseudo data” representation, we can write the discrepancy as a residual sum of squares for a regression model that is linear in both  $\mathbf{u}$  and  $\boldsymbol{\beta}$

$$d(\mathbf{u}|\mathbf{y}, \boldsymbol{\theta}, \boldsymbol{\beta}) = \left\| \begin{bmatrix} \mathbf{y} \\ 0 \end{bmatrix} - \begin{bmatrix} \mathbf{Z}\boldsymbol{\Lambda}(\boldsymbol{\theta}) & \mathbf{X} \\ \mathbf{I}_q & 0 \end{bmatrix} \begin{bmatrix} \mathbf{u} \\ \boldsymbol{\beta} \end{bmatrix} \right\|^2. \quad (11)$$

The term “pseudo data” reflects the fact that we have added  $q$  “pseudo observations” to the observed response,  $\mathbf{y}$ , and to the linear predictor,  $\boldsymbol{\gamma}(\mathbf{u}, \boldsymbol{\theta}, \boldsymbol{\beta}) = \mathbf{Z}\boldsymbol{\Lambda}(\boldsymbol{\theta})\mathbf{u} + \mathbf{X}\boldsymbol{\beta}$ , in such a way that their contribution to the overall residual sum of squares is exactly the penalty term in the discrepancy.

In the form (11) we can see that the discrepancy is a quadratic form in both  $\mathbf{u}$  and  $\boldsymbol{\beta}$ . Furthermore, because we require that  $\mathbf{X}$  has full column rank, the discrepancy is a positive-definite quadratic form in  $\mathbf{u}$  and  $\boldsymbol{\beta}$  that is minimized at  $\tilde{\mathbf{u}}(\boldsymbol{\theta})$  and  $\tilde{\boldsymbol{\beta}}(\boldsymbol{\theta})$  satisfying

$$\begin{bmatrix} \boldsymbol{\Lambda}'(\boldsymbol{\theta})\mathbf{Z}'\mathbf{Z}\boldsymbol{\Lambda}(\boldsymbol{\theta}) + \mathbf{I}_q & \boldsymbol{\Lambda}'(\boldsymbol{\theta})\mathbf{Z}'\mathbf{X} \\ \mathbf{X}'\mathbf{Z}\boldsymbol{\Lambda}(\boldsymbol{\theta}) & \mathbf{X}'\mathbf{X} \end{bmatrix} \begin{bmatrix} \tilde{\mathbf{u}}(\boldsymbol{\theta}) \\ \tilde{\boldsymbol{\beta}}(\boldsymbol{\theta}) \end{bmatrix} = \begin{bmatrix} \boldsymbol{\Lambda}'(\boldsymbol{\theta})\mathbf{Z}'\mathbf{y} \\ \mathbf{X}'\mathbf{y} \end{bmatrix} \quad (12)$$

An effective way of determining the solution to a sparse, symmetric, positive definite system of equations such as (12) is the sparse Cholesky decomposition (Davis, 2006). If  $\mathbf{A}$  is a sparse, symmetric positive definite matrix then the sparse Cholesky factor with fill-reducing permutation  $\mathbf{P}$  is the lower-triangular matrix  $\mathbf{L}$  such that

$$\mathbf{L}\mathbf{L}' = \mathbf{P}\mathbf{A}\mathbf{P}'. \quad (13)$$

(Technically, the factor  $\mathbf{L}$  is only determined up to changes in the sign of the diagonal elements. By convention we require the diagonal elements to be positive.)

The fill-reducing permutation represented by the permutation matrix  $\mathbf{P}$ , which is determined from the pattern of nonzeros in  $\mathbf{A}$  but does not depend on particular values of those nonzeros, can have a profound impact on the number of nonzeros in  $\mathbf{L}$  and hence on the speed with which  $\mathbf{L}$  can be calculated from  $\mathbf{A}$ .

In most applications of linear mixed models the matrix  $\mathbf{Z}\Lambda(\boldsymbol{\theta})$  is sparse while  $\mathbf{X}$  is dense or close to it so the permutation matrix  $\mathbf{P}$  can be restricted to the form

$$\mathbf{P} = \begin{bmatrix} \mathbf{P}_Z & \mathbf{0} \\ \mathbf{0} & \mathbf{P}_X \end{bmatrix} \quad (14)$$

without loss of efficiency. In fact, in most cases we can set  $\mathbf{P}_X = \mathbf{I}_p$  without loss of efficiency.

Let us assume that the permutation matrix is required to be of the form (14) so that we can write the Cholesky factorization for the positive definite system (12) as

$$\begin{bmatrix} \mathbf{L}_Z & \mathbf{0} \\ \mathbf{L}_{XZ} & \mathbf{L}_X \end{bmatrix} \begin{bmatrix} \mathbf{L}_Z & \mathbf{0} \\ \mathbf{L}_{XZ} & \mathbf{L}_X \end{bmatrix}' = \begin{bmatrix} \mathbf{P}_Z & \mathbf{0} \\ \mathbf{0} & \mathbf{P}_X \end{bmatrix} \begin{bmatrix} \Lambda'(\boldsymbol{\theta})\mathbf{Z}'\mathbf{Z}\Lambda(\boldsymbol{\theta}) + \mathbf{I}_q & \Lambda'(\boldsymbol{\theta})\mathbf{Z}'\mathbf{X} \\ \mathbf{X}'\mathbf{Z}\Lambda(\boldsymbol{\theta}) & \mathbf{X}'\mathbf{X} \end{bmatrix} \begin{bmatrix} \mathbf{P}_Z & \mathbf{0} \\ \mathbf{0} & \mathbf{P}_X \end{bmatrix}'. \quad (15)$$

The discrepancy can now be written in the canonical form

$$d(\mathbf{u}|\mathbf{y}, \boldsymbol{\theta}, \boldsymbol{\beta}) = \tilde{d}(\mathbf{y}, \boldsymbol{\theta}) + \left\| \begin{bmatrix} \mathbf{L}'_Z & \mathbf{L}'_{XZ} \\ \mathbf{0} & \mathbf{L}'_X \end{bmatrix} \begin{bmatrix} \mathbf{P}_Z(\mathbf{u} - \tilde{\mathbf{u}}) \\ \mathbf{P}_X(\boldsymbol{\beta} - \tilde{\boldsymbol{\beta}}) \end{bmatrix} \right\|^2 \quad (16)$$

where

$$\tilde{d}(\mathbf{y}, \boldsymbol{\theta}) = d(\tilde{\mathbf{u}}(\boldsymbol{\theta})|\mathbf{y}, \boldsymbol{\theta}, \tilde{\boldsymbol{\beta}}(\boldsymbol{\theta})) \quad (17)$$

is the minimum discrepancy, given  $\boldsymbol{\theta}$ .

### 3.2 The profiled likelihood for linear mixed models

Substituting (16) into (9) provides the unnormalized conditional density  $h(\mathbf{u}|\mathbf{y}, \boldsymbol{\theta}, \boldsymbol{\beta}, \sigma)$  on the deviance scale as

$$\begin{aligned} & -2 \log(h(\mathbf{u}|\mathbf{y}, \boldsymbol{\theta}, \boldsymbol{\beta}, \sigma)) \\ &= (n + q) \log(2\pi\sigma^2) + \frac{\tilde{d}(\mathbf{y}, \boldsymbol{\theta}) + \left\| \begin{bmatrix} \mathbf{L}'_Z & \mathbf{L}'_{XZ} \\ \mathbf{0} & \mathbf{L}'_X \end{bmatrix} \begin{bmatrix} \mathbf{P}_Z(\mathbf{u} - \tilde{\mathbf{u}}) \\ \mathbf{P}_X(\boldsymbol{\beta} - \tilde{\boldsymbol{\beta}}) \end{bmatrix} \right\|^2}{\sigma^2}. \end{aligned} \quad (18)$$



As shown in Appendix B, the integral of a quadratic form on the deviance scale, such as (18), is easily evaluated, providing the log-likelihood,  $\ell(\boldsymbol{\theta}, \boldsymbol{\beta}, \sigma | \mathbf{y})$ , as

$$\begin{aligned} -2\ell(\boldsymbol{\theta}, \boldsymbol{\beta}, \sigma | \mathbf{y}) &= -2\log(L(\boldsymbol{\theta}, \boldsymbol{\beta}, \sigma | \mathbf{y})) \\ &= n\log(2\pi\sigma^2) + \log(|\mathbf{L}\mathbf{Z}|^2) + \frac{\tilde{d}(\mathbf{y}, \boldsymbol{\theta}) + \|\mathbf{L}'_{\mathbf{X}}\mathbf{P}_{\mathbf{X}}(\boldsymbol{\beta} - \tilde{\boldsymbol{\beta}})\|^2}{\sigma^2}, \end{aligned} \quad (19)$$

from which we can see that the conditional estimate of  $\boldsymbol{\beta}$ , given  $\boldsymbol{\theta}$ , is  $\tilde{\boldsymbol{\beta}}(\boldsymbol{\theta})$  and the conditional estimate of  $\sigma$ , given  $\boldsymbol{\theta}$ , is

$$\tilde{\sigma}^2(\boldsymbol{\theta}) = \frac{\tilde{d}(\boldsymbol{\theta} | \mathbf{y})}{n}. \quad (20)$$

Substituting these conditional estimates into (19) produces the *profiled likelihood*,  $\tilde{L}(\boldsymbol{\theta} | \mathbf{y})$ , as

$$-2\tilde{\ell}(\boldsymbol{\theta} | \mathbf{y}) = \log(|\mathbf{L}\mathbf{Z}(\boldsymbol{\theta})|^2) + n \left( 1 + \log \left( \frac{2\pi\tilde{d}(\mathbf{y}, \boldsymbol{\theta})}{n} \right) \right). \quad (21)$$

The maximum likelihood estimate of  $\boldsymbol{\theta}$  can then be expressed as

$$\hat{\boldsymbol{\theta}}_L = \arg \min_{\boldsymbol{\theta}} \left( -2\tilde{\ell}(\boldsymbol{\theta} | \mathbf{y}) \right). \quad (22)$$

from which the ML estimates of  $\sigma^2$  and  $\boldsymbol{\beta}$  are evaluated as

$$\widehat{\sigma}_L^2 = \frac{\tilde{d}(\hat{\boldsymbol{\theta}}_L, \mathbf{y})}{n} \quad (23)$$

$$\hat{\boldsymbol{\beta}}_L = \tilde{\boldsymbol{\beta}}(\hat{\boldsymbol{\theta}}_L). \quad (24)$$

The important thing to note about optimizing the profiled likelihood, (21), is that it is a  $m$ -dimensional optimization problem and typically  $m$  is very small.

### 3.3 The REML criterion

In practice the so-called REML estimates of variance components are often preferred to the maximum likelihood estimates. (“REML” can be considered

to be an acronym for “restricted” or “residual” maximum likelihood, although neither term is completely accurate because these estimates do not maximize a likelihood.) We can motivate the use of the REML criterion by considering a linear regression model,

$$\mathbf{y} \sim \mathcal{N}(\mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{I}_n), \quad (25)$$

in which we typically estimate  $\sigma^2$  by

$$\widehat{\sigma}_R^2 = \frac{\|\mathbf{y} - \mathbf{X}\widehat{\boldsymbol{\beta}}\|^2}{n - p} \quad (26)$$

even though the maximum likelihood estimate of  $\sigma^2$  is

$$\widehat{\sigma}_L^2 = \frac{\|\mathbf{y} - \mathbf{X}\widehat{\boldsymbol{\beta}}\|^2}{n}. \quad (27)$$

The argument for preferring  $\widehat{\sigma}_R^2$  to  $\widehat{\sigma}_L^2$  as an estimate of  $\sigma^2$  is that the numerator in both estimates is the sum of squared residuals at  $\widehat{\boldsymbol{\beta}}$  and, although the residual vector  $\mathbf{y} - \mathbf{X}\widehat{\boldsymbol{\beta}}$  is an  $n$ -dimensional vector, the residual at  $\widehat{\boldsymbol{\theta}}$  satisfies  $p$  linearly independent constraints,  $\mathbf{X}'(\mathbf{y} - \mathbf{X}\widehat{\boldsymbol{\beta}}) = \mathbf{0}$ . That is, the residual at  $\widehat{\boldsymbol{\theta}}$  is the projection of the observed response vector,  $\mathbf{y}$ , into an  $(n - p)$ -dimensional linear subspace of the  $n$ -dimensional response space. The estimate  $\widehat{\sigma}_R^2$  takes into account the fact that  $\sigma^2$  is estimated from residuals that have only  $n - p$  *degrees of freedom*.

The REML criterion for determining parameter estimates  $\widehat{\boldsymbol{\theta}}_R$  and  $\widehat{\sigma}_R^2$  in a linear mixed model has the property that these estimates would specialize to  $\widehat{\sigma}_R^2$  from (26) for a linear regression model. Although not usually derived in this way, the REML criterion can be expressed as

$$c_R(\boldsymbol{\theta}, \boldsymbol{\sigma}|\mathbf{y}) = -2 \log \int_{\mathbb{R}^p} L(\mathbf{u}|\mathbf{y}, \boldsymbol{\theta}, \boldsymbol{\beta}, \boldsymbol{\sigma}) d\boldsymbol{\beta} \quad (28)$$

on the deviance scale. The REML estimates  $\widehat{\boldsymbol{\theta}}_R$  and  $\widehat{\sigma}_R^2$  minimize  $c_R(\boldsymbol{\theta}, \boldsymbol{\sigma}|\mathbf{y})$ .

The profiled REML criterion, a function of  $\boldsymbol{\theta}$  only, is

$$\tilde{c}_R(\boldsymbol{\theta}|\mathbf{y}) = \log(|\mathbf{L}_Z(\boldsymbol{\theta})|^2 |\mathbf{L}_X(\boldsymbol{\theta})|^2) + (n - p) \left( 1 + \log \left( \frac{2\pi \tilde{d}(\boldsymbol{\theta}|\mathbf{y})}{n - p} \right) \right) \quad (29)$$

and the REML estimate of  $\boldsymbol{\theta}$  is

$$\hat{\boldsymbol{\theta}}_R = \arg \min_{\boldsymbol{\theta}} \tilde{c}_R(\boldsymbol{\theta}, \mathbf{y}). \quad (30)$$

The REML estimate of  $\sigma^2$  is  $\hat{\sigma}_R^2 = \tilde{d}(\hat{\boldsymbol{\theta}}_R | \mathbf{y}) / (n - p)$ .

It is not entirely clear how one would define a “REML estimate” of  $\boldsymbol{\beta}$  because the REML criterion,  $c_R(\boldsymbol{\theta}, \boldsymbol{\sigma} | \mathbf{y})$ , defined in (28), does not depend on  $\boldsymbol{\beta}$ . However, it is customary (and not unreasonable) to use  $\hat{\boldsymbol{\beta}}_R = \tilde{\boldsymbol{\beta}}(\hat{\boldsymbol{\theta}}_R)$  as the REML estimate of  $\boldsymbol{\beta}$ .

Note that the profiled REML criterion can be evaluated from a sparse Cholesky decomposition like that in (15) but without the requirement that the permutation can be applied to the columns of  $\mathbf{Z}\boldsymbol{\Lambda}(\boldsymbol{\theta})$  separately from the columns of  $\mathbf{X}$ . That is, we can use a general fill-reducing permutation rather than the specific form (14) with separate permutations represented by  $\mathbf{P}_Z$  and  $\mathbf{P}_X$ . This can be useful in cases where both  $\mathbf{Z}$  and  $\mathbf{X}$  are large and sparse.

### 3.4 Summary for linear mixed models

A linear mixed model is characterized by the conditional distribution

$$(\mathcal{Y} | \mathcal{U} = \mathbf{u}) \sim \mathcal{N}(\boldsymbol{\gamma}(\mathbf{u}, \boldsymbol{\theta}, \boldsymbol{\beta}), \sigma^2 \mathbf{I}_n) \text{ where } \boldsymbol{\gamma}(\mathbf{u}, \boldsymbol{\theta}, \boldsymbol{\beta}) = \mathbf{Z}\boldsymbol{\Lambda}(\boldsymbol{\theta})\mathbf{u} + \mathbf{X}\boldsymbol{\beta} \quad (31)$$

and the unconditional distribution  $\mathcal{U} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_q)$ . The discrepancy function,

$$d(\mathbf{u} | \mathbf{y}, \boldsymbol{\theta}, \boldsymbol{\beta}) = \|\mathbf{y} - \boldsymbol{\gamma}(\mathbf{u}, \boldsymbol{\theta}, \boldsymbol{\beta})\|^2 + \|\mathbf{u}\|^2,$$

is minimized at the conditional mode,  $\tilde{\mathbf{u}}(\boldsymbol{\theta})$ , and the conditional estimate,  $\tilde{\boldsymbol{\beta}}(\boldsymbol{\theta})$ , which are the solutions to the sparse, positive-definite linear system

$$\begin{bmatrix} \boldsymbol{\Lambda}'(\boldsymbol{\theta})\mathbf{Z}'\mathbf{Z}\boldsymbol{\Lambda}(\boldsymbol{\theta}) + \mathbf{I}_q & \boldsymbol{\Lambda}'(\boldsymbol{\theta})\mathbf{Z}'\mathbf{X} \\ \mathbf{X}'\mathbf{Z}\boldsymbol{\Lambda}(\boldsymbol{\theta}) & \mathbf{X}'\mathbf{X} \end{bmatrix} \begin{bmatrix} \tilde{\mathbf{u}}(\boldsymbol{\theta}) \\ \tilde{\boldsymbol{\beta}}(\boldsymbol{\theta}) \end{bmatrix} = \begin{bmatrix} \boldsymbol{\Lambda}'(\boldsymbol{\theta})\mathbf{Z}'\mathbf{y} \\ \mathbf{X}'\mathbf{y} \end{bmatrix}.$$

In the process of solving this system we create the sparse left Cholesky factor,  $L_Z(\boldsymbol{\theta})$ , which is a lower triangular sparse matrix satisfying

$$L_Z(\boldsymbol{\theta})L_Z(\boldsymbol{\theta})' = \mathbf{P}_Z (\boldsymbol{\Lambda}'(\boldsymbol{\theta})\mathbf{Z}'\mathbf{Z}\boldsymbol{\Lambda}(\boldsymbol{\theta}) + \mathbf{I}_q) \mathbf{P}_Z'$$

where  $\mathbf{P}_Z$  is a permutation matrix representing a fill-reducing permutation formed from the pattern of nonzeros in  $\mathbf{Z}\boldsymbol{\Lambda}(\boldsymbol{\theta})$  for any  $\boldsymbol{\theta}$  not on the boundary

of the parameter region. (The values of the nonzeros depend on  $\boldsymbol{\theta}$  but the pattern doesn't.)

The profiled log-likelihood,  $\tilde{\ell}(\boldsymbol{\theta}|\mathbf{y})$ , is

$$-2\tilde{\ell}(\boldsymbol{\theta}|\mathbf{y}) = \log(|\mathbf{L}_{\mathbf{Z}}(\boldsymbol{\theta})|^2) + n \left( 1 + \log \left( \frac{2\pi\tilde{d}(\mathbf{y}, \boldsymbol{\theta})}{n} \right) \right)$$

where  $\tilde{d}(\mathbf{y}, \boldsymbol{\theta}) = d(\tilde{\mathbf{u}}(\boldsymbol{\theta})|\mathbf{y}, \tilde{\boldsymbol{\beta}}(\boldsymbol{\theta}), \boldsymbol{\theta})$ .

## 4 Generalizing the discrepancy function

Because one of the factors influencing the choice of implementation for linear mixed models is the extent to which the methods can also be applied to other mixed models, we describe several other classes of mixed models before discussing the implementation details for linear mixed models. At the core of our methods for determining the maximum likelihood estimates (MLEs) of the parameters in the mixed model are methods for minimizing the discrepancy function with respect to the coefficients  $\mathbf{u}$  and  $\boldsymbol{\beta}$  in the linear predictor  $\gamma(\mathbf{u}, \boldsymbol{\theta}, \boldsymbol{\beta})$ .

In this section we describe the general form of the discrepancy function that we will use and a penalized iteratively reweighted least squares (PIRLS) algorithm for determining the conditional modes  $\tilde{\mathbf{u}}(\boldsymbol{\theta})$  and  $\tilde{\boldsymbol{\beta}}(\boldsymbol{\theta})$ . We then describe several types of mixed models and the form of the discrepancy function for each.

### 4.1 A weighted residual sum of squares

As shown in §3.1, the discrepancy function for a linear mixed model has the form of a penalized residual sum of squares from a linear model (11). In this section we generalize that definition to

$$d(\mathbf{u}|\mathbf{y}, \boldsymbol{\theta}, \boldsymbol{\beta}) = \|\mathbf{W}^{1/2}(\boldsymbol{\mu}) [\mathbf{y} - \boldsymbol{\mu}_{\mathcal{Y}|\mathcal{U}}(\mathbf{u}, \boldsymbol{\theta}, \boldsymbol{\beta})]\|^2 + \|\mathbf{0} - \mathbf{u}\|^2. \quad (32)$$

where  $\mathbf{W}$  is an  $n \times n$  diagonal matrix, called the *weights matrix*, with positive diagonal elements and  $\mathbf{W}^{1/2}$  is the diagonal matrix with the square roots of the weights on the diagonal. The  $i$ th weight is inversely proportional to the conditional variances of  $(\mathcal{Y}|\mathcal{U} = \mathbf{u})$  and may depend on the conditional mean,  $\boldsymbol{\mu}_{\mathcal{Y}|\mathcal{U}}$ .

We allow the conditional mean to be a nonlinear function of the linear predictor, but with certain restrictions. We require that the mapping from  $\mathbf{u}$  to  $\boldsymbol{\mu}_{\mathcal{Y}|\mathcal{U}=\mathbf{u}}$  be expressed as

$$\mathbf{u} \rightarrow \boldsymbol{\gamma} \rightarrow \boldsymbol{\eta} \rightarrow \boldsymbol{\mu} \quad (33)$$

where  $\boldsymbol{\gamma} = \mathbf{Z}\boldsymbol{\Lambda}(\boldsymbol{\theta})\mathbf{u} + \mathbf{X}\boldsymbol{\theta}$  is an  $ns$ -dimensional vector ( $s > 0$ ) while  $\boldsymbol{\eta}$  and  $\boldsymbol{\mu}$  are  $n$ -dimensional vectors.

The map  $\boldsymbol{\eta} \rightarrow \boldsymbol{\mu}$  has the property that  $\mu_i$  depends only on  $\eta_i$ ,  $i = 1, \dots, n$ . The map  $\boldsymbol{\gamma} \rightarrow \boldsymbol{\eta}$  has a similar property in that, if we write  $\boldsymbol{\gamma}$  as an  $n \times s$  matrix  $\boldsymbol{\Gamma}$  such that

$$\boldsymbol{\gamma} = \text{vec } \boldsymbol{\Gamma} \quad (34)$$

(i.e. concatenating the columns of  $\boldsymbol{\Gamma}$  produces  $\boldsymbol{\gamma}$ ) then  $\eta_i$  depends only on the  $i$ th row of  $\boldsymbol{\Gamma}$ ,  $i = 1, \dots, n$ . Thus the Jacobian matrix  $\frac{d\boldsymbol{\mu}}{d\boldsymbol{\eta}}$  is an  $n \times n$  diagonal matrix and the Jacobian matrix  $\frac{d\boldsymbol{\eta}}{d\boldsymbol{\gamma}}$  is the horizontal concatenation of  $s$  diagonal  $n \times n$  matrices.

For historical reasons, the function that maps  $\eta_i$  to  $\mu_i$  is called the *inverse link* function and is written  $\mu = g^{-1}(\eta)$ . The *link function*, naturally, is  $\eta = g(\mu)$ . When applied component-wise to vectors  $\boldsymbol{\mu}$  or  $\boldsymbol{\eta}$  we write these as  $\boldsymbol{\eta} = \mathbf{g}(\boldsymbol{\mu})$  and  $\boldsymbol{\mu} = \mathbf{g}^{-1}(\boldsymbol{\eta})$ .

Recall that the conditional distribution,  $(\mathcal{Y}_i|\mathcal{U} = \mathbf{u})$ , is required to be independent of  $(\mathcal{Y}_j|\mathcal{U} = \mathbf{u})$  for  $i, j = 1, \dots, n$ ,  $i \neq j$  and that all the component conditional distributions must be of the same form and differ only according to the value of the conditional mean.

Depending on the family of the conditional distributions, the allowable values of the  $\mu_i$  may be in a restricted range. For example, if the conditional distributions are Bernoulli then  $0 \leq \mu_i \leq 1$ ,  $i = 1, \dots, n$ . If the conditional distributions are Poisson then  $0 \leq \mu_i$ ,  $i = 1, \dots, n$ . A characteristic of the link function,  $g$ , is that it must map the restricted range to an unrestricted range. That is, a link function for the Bernoulli distribution must map  $[0, 1]$  to  $[-\infty, \infty]$  and must be invertible within the range.

The mapping from  $\boldsymbol{\gamma}$  to  $\boldsymbol{\eta}$  is defined by a function  $m : \mathbb{R}^s \rightarrow \mathbb{R}$ , called the *nonlinear model* function, such that  $\eta_i = m(\boldsymbol{\gamma}_i)$ ,  $i = 1, \dots, n$  where  $\boldsymbol{\gamma}_i$  is the  $i$ th row of  $\boldsymbol{\Gamma}$ . The vector-valued function is  $\boldsymbol{\eta} = \mathbf{m}(\boldsymbol{\gamma})$ .

Determining the conditional modes,  $\tilde{\mathbf{u}}(\mathbf{y}|\boldsymbol{\theta})$ , and  $\tilde{\boldsymbol{\beta}}(\mathbf{y}|\boldsymbol{\theta})$ , that jointly minimize the discrepancy,

$$\begin{bmatrix} \tilde{\mathbf{u}}(\mathbf{y}|\boldsymbol{\theta}) \\ \tilde{\boldsymbol{\beta}}(\mathbf{y}|\boldsymbol{\theta}) \end{bmatrix} = \arg \min_{\mathbf{u}, \boldsymbol{\beta}} [(\mathbf{y} - \boldsymbol{\mu})' \mathbf{W}(\mathbf{y} - \boldsymbol{\mu}) + \|\mathbf{u}\|^2] \quad (35)$$

becomes a weighted, nonlinear least squares problem except that the weights,  $\mathbf{W}$ , can depend on  $\boldsymbol{\mu}$  and, hence, on  $\mathbf{u}$  and  $\boldsymbol{\beta}$ .

In describing an algorithm for linear mixed models we called  $\tilde{\boldsymbol{\beta}}(\boldsymbol{\theta})$  the *conditional estimate*. That name reflects that fact that this is the maximum likelihood estimate of  $\boldsymbol{\beta}$  for that particular value of  $\boldsymbol{\theta}$ . Once we have determined the MLE,  $\hat{\boldsymbol{\theta}}_L$  of  $\boldsymbol{\theta}$ , we have a “plug-in” estimator,  $\hat{\boldsymbol{\beta}}_L = \tilde{\boldsymbol{\beta}}(\boldsymbol{\theta})$  for  $\boldsymbol{\beta}$ .

This property does not carry over exactly to other forms of mixed models. The values  $\tilde{\mathbf{u}}(\boldsymbol{\theta})$  and  $\tilde{\boldsymbol{\beta}}(\boldsymbol{\theta})$  are conditional modes in the sense that they are the coefficients in  $\boldsymbol{\gamma}$  that jointly maximize the unscaled conditional density  $h(\mathbf{u}|\mathbf{y}, \boldsymbol{\theta}, \boldsymbol{\beta}, \sigma)$ . Here we are using the adjective “conditional” more in the sense of conditioning on  $\mathbf{Y} = \mathbf{y}$  than in the sense of conditioning on  $\boldsymbol{\theta}$ , although these values are determined for a fixed value of  $\boldsymbol{\theta}$ .

## 4.2 The PIRLS algorithm for $\tilde{\mathbf{u}}$ and $\tilde{\boldsymbol{\beta}}$

The penalized, iteratively reweighted, least squares (PIRLS) algorithm to determine  $\tilde{\mathbf{u}}(\boldsymbol{\theta})$  and  $\tilde{\boldsymbol{\beta}}(\boldsymbol{\theta})$  is a form of the Fisher scoring algorithm. We fix the weights matrix,  $\mathbf{W}$ , and use penalized, weighted, nonlinear least squares to minimize the penalized, weighted residual sum of squares conditional on these weights. Then we update the weights to those determined by the current value of  $\boldsymbol{\mu}$  and iterate.

To describe this algorithm in more detail we will use parenthesized superscripts to denote the iteration number. Thus  $\mathbf{u}^{(0)}$  and  $\boldsymbol{\beta}^{(0)}$  are the initial values of these parameters, while  $\mathbf{u}^{(i)}$  and  $\boldsymbol{\beta}^{(i)}$  are the values at the  $i$ th iteration. Similarly  $\boldsymbol{\gamma}^{(i)} = \mathbf{Z}\boldsymbol{\Lambda}(\boldsymbol{\theta})\mathbf{u}^{(i)} + \mathbf{X}\boldsymbol{\beta}^{(i)}$ ,  $\boldsymbol{\eta}^{(i)} = \mathbf{m}(\boldsymbol{\gamma}^{(i)})$  and  $\boldsymbol{\mu}^{(i)} = \mathbf{g}^{-1}(\boldsymbol{\eta}^{(i)})$ .

We use a penalized version of the Gauss-Newton algorithm (Bates and Watts, 1988, ch. 2) for which we define the weighted Jacobian matrices

$$\mathbf{U}^{(i)} = \mathbf{W}^{1/2} \left. \frac{d\boldsymbol{\mu}}{d\mathbf{u}'} \right|_{\mathbf{u}=\mathbf{u}^{(i)}, \boldsymbol{\beta}=\boldsymbol{\beta}^{(i)}} = \mathbf{W}^{1/2} \left. \frac{d\boldsymbol{\mu}}{d\boldsymbol{\eta}'} \right|_{\boldsymbol{\eta}^{(i)}} \left. \frac{d\boldsymbol{\eta}}{d\boldsymbol{\gamma}'} \right|_{\boldsymbol{\gamma}^{(i)}} \mathbf{Z}\boldsymbol{\Lambda}(\boldsymbol{\theta}) \quad (36)$$

$$\mathbf{V}^{(i)} = \mathbf{W}^{1/2} \left. \frac{d\boldsymbol{\mu}}{d\boldsymbol{\beta}'} \right|_{\mathbf{u}=\mathbf{u}^{(i)}, \boldsymbol{\beta}=\boldsymbol{\beta}^{(i)}} = \mathbf{W}^{1/2} \left. \frac{d\boldsymbol{\mu}}{d\boldsymbol{\eta}'} \right|_{\boldsymbol{\eta}^{(i)}} \left. \frac{d\boldsymbol{\eta}}{d\boldsymbol{\gamma}'} \right|_{\boldsymbol{\gamma}^{(i)}} \mathbf{X} \quad (37)$$

of dimension  $n \times q$  and  $n \times p$ , respectively. The increments at the  $i$ th iteration,  $\boldsymbol{\delta}_u^{(i)}$  and  $\boldsymbol{\delta}_\beta^{(i)}$ , are the solutions to

$$\begin{bmatrix} \mathbf{U}^{(i)'}\mathbf{U}^{(i)} + \mathbf{I}_q & \mathbf{U}^{(i)'}\mathbf{V}^{(i)} \\ \mathbf{V}^{(i)'}\mathbf{U}^{(i)} & \mathbf{V}^{(i)'}\mathbf{V}^{(i)} \end{bmatrix} \begin{bmatrix} \boldsymbol{\delta}_u^{(i)} \\ \boldsymbol{\delta}_\beta^{(i)} \end{bmatrix} = \begin{bmatrix} \mathbf{U}^{(i)'}\mathbf{W}^{1/2}(\mathbf{y} - \boldsymbol{\mu}^{(i)}) - \mathbf{u}^{(i)} \\ \mathbf{U}^{(i)'}\mathbf{W}^{1/2}(\mathbf{y} - \boldsymbol{\mu}^{(i)}) \end{bmatrix} \quad (38)$$

providing the updated parameter values

$$\begin{bmatrix} \mathbf{u}^{(i+1)} \\ \boldsymbol{\beta}^{(i+1)} \end{bmatrix} = \begin{bmatrix} \mathbf{u}^{(i)} \\ \boldsymbol{\beta}^{(i)} \end{bmatrix} + \lambda \begin{bmatrix} \boldsymbol{\delta}_u^{(i)} \\ \boldsymbol{\delta}_\beta^{(i)} \end{bmatrix} \quad (39)$$

where  $\lambda > 0$  is a step factor chosen to ensure that

$$(\mathbf{y} - \boldsymbol{\mu}^{(i+1)})' \mathbf{W} (\mathbf{y} - \boldsymbol{\mu}^{(i+1)}) + \|\mathbf{u}^{(i+1)}\|^2 < (\mathbf{y} - \boldsymbol{\mu}^{(i)})' \mathbf{W} (\mathbf{y} - \boldsymbol{\mu}^{(i)}) + \|\mathbf{u}^{(i)}\|^2. \quad (40)$$

In the process of solving for the increments we form the sparse, lower triangular, Cholesky factor,  $\mathbf{L}^{(i)}$ , satisfying

$$\mathbf{L}^{(i)} \mathbf{L}^{(i)'} = \mathbf{P}_Z \left( \mathbf{U}^{(i)'} \mathbf{U}^{(i)} + \mathbf{I}_n \right) \mathbf{P}_Z'. \quad (41)$$

After each successful iteration, determining new values of the coefficients,  $\mathbf{u}^{(i+1)}$  and  $\boldsymbol{\beta}^{(i+1)}$ , that reduce the penalized, weighted residual sum of squares, we update the weights matrix to  $\mathbf{W}(\boldsymbol{\mu}^{(i+1)})$  and the weighted Jacobians,  $\mathbf{U}^{(i+1)}$  and  $\mathbf{V}^{(i+1)}$ , then iterate. Convergence is determined according to the orthogonality convergence criterion (Bates and Watts, 1988, ch. 2), suitably adjusted for the weights matrix and the penalty.

### 4.3 Weighted linear mixed models

One of the simplest generalizations of linear mixed models is a weighted linear mixed model where  $s = 1$ , the link function,  $g$ , and the nonlinear model function,  $m$ , are both the identity, the weights matrix,  $\mathbf{W}$ , is constant and the conditional distribution family is Gaussian. That is, the conditional distribution can be written

$$(\mathbf{y} | \mathbf{u} = \mathbf{u}) \sim \mathcal{N}(\gamma(\mathbf{u}, \boldsymbol{\theta}, \boldsymbol{\beta}), \sigma^2 \mathbf{W}^{-1}) \quad (42)$$

with discrepancy function

$$d(\mathbf{u} | \mathbf{y}, \boldsymbol{\theta}, \boldsymbol{\beta}) = \|\mathbf{W}^{1/2}(\mathbf{y} - \mathbf{Z}\boldsymbol{\Lambda}(\boldsymbol{\theta})\mathbf{u} - \mathbf{X}\boldsymbol{\theta})\|^2 + \|\mathbf{u}\|^2. \quad (43)$$

The conditional mode,  $\tilde{\mathbf{u}}(\boldsymbol{\theta})$ , and the conditional estimate,  $\tilde{\boldsymbol{\beta}}(\boldsymbol{\theta})$ , are the solutions to

$$\begin{bmatrix} \boldsymbol{\Lambda}'(\boldsymbol{\theta}) \mathbf{Z}' \mathbf{W} \mathbf{Z} \boldsymbol{\Lambda}(\boldsymbol{\theta}) + \mathbf{I}_q & \boldsymbol{\Lambda}'(\boldsymbol{\theta}) \mathbf{Z}' \mathbf{W} \mathbf{X} \\ \mathbf{X}' \mathbf{W} \mathbf{Z} \boldsymbol{\Lambda}(\boldsymbol{\theta}) & \mathbf{X}' \mathbf{W} \mathbf{X} \end{bmatrix} \begin{bmatrix} \tilde{\mathbf{u}}(\boldsymbol{\theta}) \\ \tilde{\boldsymbol{\beta}}(\boldsymbol{\theta}) \end{bmatrix} = \begin{bmatrix} \boldsymbol{\Lambda}'(\boldsymbol{\theta}) \mathbf{Z}' \mathbf{W} \mathbf{y} \\ \mathbf{X}' \mathbf{W} \mathbf{y} \end{bmatrix}, \quad (44)$$

which can be solved directly, and the Cholesky factor,  $\mathbf{L}_Z(\boldsymbol{\theta})$ , satisfies

$$\mathbf{L}_Z(\boldsymbol{\theta})\mathbf{L}_Z(\boldsymbol{\theta})' = \mathbf{P}_Z (\boldsymbol{\Lambda}'(\boldsymbol{\theta})\mathbf{Z}'\mathbf{W}\mathbf{Z}\boldsymbol{\Lambda}(\boldsymbol{\theta}) + \mathbf{I}_q) \mathbf{P}_Z'. \quad (45)$$

The profiled log-likelihood,  $\tilde{\ell}(\boldsymbol{\theta}|\mathbf{y})$ , is

$$-2\tilde{\ell}(\boldsymbol{\theta}|\mathbf{y}) = \log \left( \frac{|\mathbf{L}_Z(\boldsymbol{\theta})|^2}{|\mathbf{W}|} \right) + n \left( 1 + \log \left( \frac{2\pi\tilde{d}(\mathbf{y}, \boldsymbol{\theta})}{n} \right) \right). \quad (46)$$

If the matrix  $\mathbf{W}$  is fixed then we can ignore the term  $|\mathbf{W}|$  in (46) when determining the MLE,  $\hat{\boldsymbol{\theta}}_L$ . However, in some models, we use a parameterized weight matrix,  $\mathbf{W}(\boldsymbol{\phi})$ , and wish to determine the MLEs,  $\hat{\boldsymbol{\phi}}_L$  and  $\hat{\boldsymbol{\theta}}_L$  simultaneously. In these cases we must include the term involving  $|\mathbf{W}(\boldsymbol{\phi})|$  when evaluating the profiled log-likelihood.

Note that we must define the parameterization of  $\mathbf{W}(\boldsymbol{\phi})$  such that  $\sigma^2$  and  $\boldsymbol{\phi}$  are not a redundant parameterization of  $\sigma^2\mathbf{W}(\boldsymbol{\phi})$ . For example, we could require that the first diagonal element of  $\mathbf{W}$  be unity.

#### 4.4 Nonlinear mixed models

In an unweighted, nonlinear mixed model the conditional distribution is Gaussian, the link,  $g$ , is the identity and the weights matrix,  $\mathbf{W} = \mathbf{I}_n$ . That is,

$$(\mathbf{y}|\mathbf{u} = \mathbf{u}) \sim \mathcal{N}(\mathbf{m}(\boldsymbol{\gamma}), \sigma^2 \mathbf{I}_n) \quad (47)$$

with discrepancy function

$$d(\mathbf{u}|\mathbf{y}, \boldsymbol{\theta}, \boldsymbol{\beta}) = \|\mathbf{y} - \boldsymbol{\mu}\|^2 + \|\mathbf{u}\|^2. \quad (48)$$

For a given value of  $\boldsymbol{\theta}$  we determine the conditional modes,  $\tilde{\mathbf{u}}(\boldsymbol{\theta})$  and  $\tilde{\boldsymbol{\beta}}(\boldsymbol{\theta})$ , as the solution to the penalized nonlinear least squares problem

$$\begin{bmatrix} \tilde{\mathbf{u}}(\boldsymbol{\theta}) \\ \tilde{\boldsymbol{\beta}}(\boldsymbol{\theta}) \end{bmatrix} = \arg \min_{\mathbf{u}, \boldsymbol{\beta}} d(\mathbf{u}|\mathbf{y}, \boldsymbol{\theta}, \boldsymbol{\beta}) \quad (49)$$

and we write the minimum discrepancy, given  $\mathbf{y}$  and  $\boldsymbol{\theta}$ , as

$$\tilde{d}(\mathbf{y}, \boldsymbol{\theta}) = d(\tilde{\mathbf{u}}(\boldsymbol{\theta})|\mathbf{y}, \boldsymbol{\theta}, \tilde{\boldsymbol{\beta}}(\boldsymbol{\theta})). \quad (50)$$



Let  $\tilde{\mathbf{L}}_Z(\boldsymbol{\theta})$  and  $\tilde{\mathbf{L}}_X(\boldsymbol{\theta})$  be the Cholesky factors at  $\boldsymbol{\theta}$ ,  $\tilde{\boldsymbol{\beta}}(\boldsymbol{\theta})$  and  $\tilde{\mathbf{u}}(\boldsymbol{\theta})$ . Then the *Laplace approximation* to the log-likelihood is

$$-2\ell_P(\boldsymbol{\theta}, \boldsymbol{\beta}, \sigma|\mathbf{y}) \approx n \log(2\pi\sigma^2) + \log(|\tilde{\mathbf{L}}_Z|^2) + \frac{\tilde{d}(\mathbf{y}, \boldsymbol{\theta}) + \left\| \tilde{\mathbf{L}}'_X(\boldsymbol{\beta} - \tilde{\boldsymbol{\beta}}) \right\|^2}{\sigma^2}, \quad (51)$$

producing the approximate profiled log-likelihood,  $\tilde{\ell}_P(\boldsymbol{\theta}|\mathbf{y})$ ,

$$-2\tilde{\ell}_P(\boldsymbol{\theta}|\mathbf{y}) \approx \log(|\tilde{\mathbf{L}}_Z|^2) + n \left( 1 + \log(2\pi\tilde{d}(\mathbf{y}, \boldsymbol{\theta})/n) \right). \quad (52)$$

#### 4.4.1 Nonlinear mixed model summary

In a nonlinear mixed model we determine the parameter estimate,  $\hat{\boldsymbol{\theta}}_P$ , from the Laplace approximation to the log-likelihood as

$$\hat{\boldsymbol{\theta}}_P = \arg \max_{\boldsymbol{\theta}} \tilde{\ell}_P(\boldsymbol{\theta}|\mathbf{y}) = \arg \min_{\boldsymbol{\theta}} \log(|\tilde{\mathbf{L}}_Z|^2) + n \left( 1 + \log(2\pi\tilde{d}(\mathbf{y}, \boldsymbol{\theta})/n) \right). \quad (53)$$

Each evaluation of  $\tilde{\ell}_P(\boldsymbol{\theta}|\mathbf{y})$  requires a solving the penalized nonlinear least squares problem (49) simultaneously with respect to both sets of coefficients,  $\mathbf{u}$  and  $\boldsymbol{\beta}$ , in the linear predictor,  $\boldsymbol{\gamma}$ .

For a weighted nonlinear mixed model with fixed weights,  $\mathbf{W}$ , we replace the unweighted discrepancy function  $d(\mathbf{u}|\mathbf{y}, \boldsymbol{\theta}, \boldsymbol{\beta})$  with the weighted discrepancy function,

## 5 Details of the implementation

### 5.1 Implementation details for linear mixed models

The crucial step in implementing algorithms for determining ML or REML estimates of the parameters in a linear mixed model is evaluating the factorization (15) for any  $\boldsymbol{\theta}$  satisfying  $\boldsymbol{\theta}_L \leq \boldsymbol{\theta} \leq \boldsymbol{\theta}_U$ . We will assume that  $\mathbf{Z}$  is sparse as is  $\mathbf{Z}\boldsymbol{\Lambda}(\boldsymbol{\theta})$ .

When  $\mathbf{X}$  is not sparse we will use the factorization (15) setting  $\mathbf{P}_X = \mathbf{I}_p$  and storing  $\mathbf{L}_{XZ}$  and  $\mathbf{L}_X$  as dense matrices. The permutation matrix  $\mathbf{P}_Z$  is determined from the pattern of non-zeros in  $\mathbf{Z}\boldsymbol{\Lambda}(\boldsymbol{\theta})$  which does not depend on  $\boldsymbol{\theta}$ , as long as  $\boldsymbol{\theta}$  is not on the boundary. In fact, in most cases the pattern of non-zeros in  $\mathbf{Z}\boldsymbol{\Lambda}(\boldsymbol{\theta})$  is the same as the pattern of non-zeros in  $\mathbf{Z}$ . For many

models, in particular models with scalar random effects (described later), the matrix  $\Lambda(\boldsymbol{\theta})$  is diagonal.

Given a value of  $\boldsymbol{\theta}$  we determine the Cholesky factor  $\mathbf{L}_Z$  satisfying

$$\mathbf{L}_Z \mathbf{L}'_Z = \mathbf{P}_Z (\Lambda'(\boldsymbol{\theta}) \mathbf{Z}' \mathbf{Z} \Lambda(\boldsymbol{\theta}) + \mathbf{I}_q) \mathbf{P}'_Z. \quad (54)$$

The CHOLMOD package allows for  $\mathbf{L}_Z$  to be calculated directly from  $\Lambda'(\boldsymbol{\theta}) \mathbf{Z}'$  or from  $\Lambda'(\boldsymbol{\theta}) \mathbf{Z}' \mathbf{Z} \Lambda(\boldsymbol{\theta})$ . The choice in implementation is whether to store  $\mathbf{Z}'$  and update it to  $\Lambda'(\boldsymbol{\theta}) \mathbf{Z}'$  or to store  $\mathbf{Z}' \mathbf{Z}$  and use it to form  $\Lambda'(\boldsymbol{\theta}) \mathbf{Z}' \mathbf{Z} \Lambda(\boldsymbol{\theta})$  at each evaluation.

In the lme4 package we store  $\mathbf{Z}'$  and use it to form  $\Lambda'(\boldsymbol{\theta}) \mathbf{Z}'$  from which  $\mathbf{L}_Z$  is evaluated. There are two reasons for this choice. First, the calculations for the more general forms of mixed models cannot be reduced to calculations involving  $\mathbf{Z}' \mathbf{Z}$  and by expressing these calculations in terms of  $\Lambda(\boldsymbol{\theta}) \mathbf{Z}'$  for linear mixed models we can reuse the code for the more general models. Second, the calculation of  $\Lambda(\boldsymbol{\theta})' (\mathbf{Z}' \mathbf{Z}) \Lambda(\boldsymbol{\theta})$  from  $\mathbf{Z}' \mathbf{Z}$  is complicated compared to the calculation of  $\Lambda(\boldsymbol{\theta})' \mathbf{Z}'$  from  $\mathbf{Z}'$ .

This choice is disadvantageous when  $n \gg q$  because  $\mathbf{Z}'$  is much larger than  $\mathbf{Z}' \mathbf{Z}$ , even when they are stored as sparse matrices. Evaluation of  $\mathbf{L}_Z$  directly from  $\mathbf{Z}'$  requires more storage and more calculation than evaluating  $\mathbf{L}_Z$  from  $\mathbf{Z}' \mathbf{Z}$ .

Next we evaluate  $\mathbf{L}'_{XZ}$  as the solution to

$$\mathbf{L}_Z \mathbf{L}'_{XZ} = \mathbf{P}_Z \Lambda'(\boldsymbol{\theta}) \mathbf{Z}' \mathbf{X}. \quad (55)$$

Again we have the choice of calculating and storing  $\mathbf{Z}' \mathbf{X}$  or storing  $\mathbf{X}$  and using it to reevaluate  $\mathbf{Z}' \mathbf{X}$ . In the lme4 package we store  $\mathbf{X}$ , because the calculations for the more general models cannot be expressed in terms of  $\mathbf{Z}' \mathbf{X}$ .

Finally  $\mathbf{L}_X$  is evaluated as the (dense) solution to

$$\mathbf{L}_X \mathbf{L}'_X = \mathbf{X}' \mathbf{X} - \mathbf{L}_{XZ} \mathbf{L}_{XZ}. \quad (56)$$

from which  $\tilde{\boldsymbol{\beta}}$  can be determined as the solution to dense system

$$\mathbf{L}_X \mathbf{L}_X \tilde{\boldsymbol{\beta}} = \mathbf{X}' \mathbf{y} \quad (57)$$

and  $\tilde{\mathbf{u}}$  as the solution to the sparse system

$$\mathbf{L}_Z \mathbf{L}_Z \tilde{\mathbf{u}} = \Lambda' \mathbf{Z}' \mathbf{y} \quad (58)$$

For many models, in particular models with scalar random effects, which are described later, the matrix  $\Lambda(\boldsymbol{\theta})$  is diagonal. For such a model, if both  $\mathbf{Z}$  and  $\mathbf{X}$  are sparse and we plan to use the REML criterion then we create and store

$$\mathbf{A} = \begin{bmatrix} \mathbf{Z}'\mathbf{Z} & \mathbf{Z}'\mathbf{X} \\ \mathbf{X}'\mathbf{Z} & \mathbf{X}'\mathbf{X} \end{bmatrix} \quad \text{and} \quad \mathbf{c} = \begin{bmatrix} \mathbf{Z}'\mathbf{y} \\ \mathbf{X}'\mathbf{y} \end{bmatrix} \quad (59)$$

and determine a fill-reducing permutation,  $\mathbf{P}$ , for  $\mathbf{A}$ . Given a value of  $\boldsymbol{\theta}$  we create the factorization

$$\mathbf{L}(\boldsymbol{\theta})\mathbf{L}(\boldsymbol{\theta})' = \mathbf{P} \left( \begin{bmatrix} \Lambda(\boldsymbol{\theta}) & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_p \end{bmatrix} \mathbf{A} \begin{bmatrix} \Lambda(\boldsymbol{\theta}) & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_p \end{bmatrix} + \begin{bmatrix} \mathbf{I}_q & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \right) \mathbf{P}' \quad (60)$$

solve for  $\tilde{\mathbf{u}}(\boldsymbol{\theta})$  and  $\tilde{\boldsymbol{\beta}}(\boldsymbol{\theta})$  in

$$\mathbf{L}\mathbf{L}'\mathbf{P} \begin{bmatrix} \tilde{\mathbf{u}}(\boldsymbol{\theta}) \\ \tilde{\boldsymbol{\beta}}(\boldsymbol{\theta}) \end{bmatrix} = \mathbf{P} \begin{bmatrix} \Lambda(\boldsymbol{\theta}) & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_p \end{bmatrix} \mathbf{c} \quad (61)$$

then evaluate  $\tilde{d}(\mathbf{y}|\boldsymbol{\theta})$  and the profiled REML criterion as

$$\tilde{d}_R(\boldsymbol{\theta}|\mathbf{y}) = \log(|\mathbf{L}(\boldsymbol{\theta})|^2) + (n - p) \left( 1 + \log \left( \frac{2\pi\tilde{d}(\mathbf{y}|\boldsymbol{\theta})}{n - p} \right) \right). \quad (62)$$

## References

- Douglas M. Bates and Donald G. Watts. *Nonlinear Regression Analysis and Its Applications*. Wiley, 1988.
- Tim Davis. CHOLMOD: sparse supernodal Cholesky factorization and update/downdate. <http://www.cise.ufl.edu/research/sparse/cholmod>, 2005.
- Timothy A. Davis. *Direct Methods for Sparse Linear Systems*. Fundamentals of Algorithms. SIAM, 2006.

## A Notation

### A.1 Random variables in the model

$\mathbf{B}$  Random-effects vector of dimension  $q$ ,  $\mathbf{B} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{V}(\boldsymbol{\theta}) \mathbf{V}(\boldsymbol{\theta})')$ .

$\mathbf{U}$  “Spherical” random-effects vector of dimension  $q$ ,  $\mathbf{U} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_q)$ ,  $\mathbf{B} = \mathbf{V}(\boldsymbol{\theta})\mathbf{U}$ .

$\mathbf{Y}$  Response vector of dimension  $n$ .

## A.2 Parameters of the model

$\boldsymbol{\beta}$  Fixed-effects parameters (dimension  $p$ ).

$\boldsymbol{\theta}$  Parameters determining the left factor,  $\mathbf{\Lambda}(\boldsymbol{\theta})$  of the relative covariance matrix of  $\mathbf{B}$  (dimension  $m$ ).

$\sigma$  the common scale parameter - not used in some generalized linear mixed models and generalized nonlinear mixed models.

## A.3 Dimensions

$m$  dimension of the parameter  $\boldsymbol{\theta}$ .

$n$  dimension of the response vector,  $\mathbf{y}$ , and the random variable,  $\mathbf{Y}$ .

$p$  dimension of the fixed-effects parameter,  $\boldsymbol{\beta}$ .

$q$  dimension of the random effects,  $\mathbf{B}$  or  $\mathbf{U}$ .

$s$  dimension of the parameter vector,  $\boldsymbol{\phi}$ , in the nonlinear model function.

## A.4 Matrices

$\mathbf{L}$  Left Cholesky factor of a positive-definite symmetric matrix.  $\mathbf{L}_Z$  is  $q \times q$ ;  $\mathbf{L}_X$  is  $p \times p$ .

$\mathbf{P}$  Fill-reducing permutation for the random effects model matrix. (Size  $q \times q$ .)

$\mathbf{V}$  Left factor of the relative covariance matrix of the random effects. (Size  $q \times q$ .)

$\mathbf{X}$  Model matrix for the fixed-effects parameters,  $\boldsymbol{\beta}$ . (Size  $(ns) \times p$ .)

$\mathbf{Z}$  Model matrix for the random effects. (Size  $(ns) \times q$ .)

## B Integrating a quadratic deviance expression

In (6) we defined the likelihood of the parameters given the observed data as

$$L(\boldsymbol{\theta}, \boldsymbol{\beta}, \sigma | \mathbf{y}) = \int_{\mathbb{R}^q} h(\mathbf{u} | \mathbf{y}, \boldsymbol{\theta}, \boldsymbol{\beta}, \sigma) d\mathbf{u}.$$

which is often alarmingly described as “an intractable integral”. In point of fact, this integral can be evaluated exactly in the case of a linear mixed model and can be approximated quite accurately for other forms of mixed models.