

Computational methods for mixed models

Douglas Bates
Department of Statistics
University of Wisconsin – Madison

September 9, 2008

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 a particular form of mixed model, called a linear mixed model, and also the general form of the mixed models that can be represented in the `lme4` package.

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,

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

The scalar, σ , 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 σ^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 $\mathbf{\Lambda}(\boldsymbol{\theta})$, 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 θ_{iL} and θ_{iU} , $i = 1, \dots, m$. We shall write the set of such constraints as $\boldsymbol{\theta}_L \leq \boldsymbol{\theta} \leq \boldsymbol{\theta}_U$. The matrix $\mathbf{\Lambda}(\boldsymbol{\theta})$ 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 the following conditions:

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. Frequently, the number of rows in \mathbf{Z} and \mathbf{X} is n , the dimension of \mathbf{y} , but not always.

2. The scalar distributions, $(\mathcal{Y}_i|\mathbf{B} = \mathbf{b}), i = 1, \dots, n$, are independent. We say that the components of \mathbf{Y} are *conditionally independent* given \mathbf{B} . Furthermore, the scalar distributions must all have the same form. They differ only in the value of the conditional means, $\mu_i(\mathbf{b})$.
3. The scalar distributions, $(\mathcal{Y}_i|\mathbf{B} = \mathbf{b}), i = 1, \dots, n$, are completely determined by the conditional mean, $\mu_{\mathbf{Y}|\mathbf{B}}(\mathbf{b})$ and, at most, one additional parameter, σ , which is the common scale parameter.

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

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

where \mathbf{I}_n denotes the identity matrix of size n . In this case the conditional mean, $\mu_{\mathbf{Y}|\mathbf{B}}(\mathbf{b})$, is exactly the linear predictor, $\mathbf{Zb} + \mathbf{X}\boldsymbol{\beta}$, 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 require that we 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, defined by the 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 σ , given the observed data \mathbf{y} . The *maximum likelihood estimates* of these parameters are the values that maximize L .

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, ($\mathcal{Y}|\mathcal{U} = \mathbf{u}$), and the unconditional distribution, \mathcal{U} , are spherical Gaussian distributions and that the conditional mean, $\boldsymbol{\mu}_{\mathcal{Y}|\mathcal{U}}(\mathbf{u})$, is the linear predictor, $\boldsymbol{\gamma}(\mathbf{u})$. Because all the distributions determining the model are continuous distributions, we consider their densities. It is convenient to express these densities on the scale of the logarithm of the density or, as

is common in many statistical expressions, on the *deviance scale* (negative twice the logarithm of the density),

$$\begin{aligned}
-2\log(f_{\mathbf{u}}(\mathbf{u})) &= q\log(2\pi\sigma^2) + \frac{\|\mathbf{u}\|^2}{\sigma^2} \\
-2\log(f_{\mathbf{y}|\mathbf{u}}(\mathbf{y}|\mathbf{u})) &= n\log(2\pi\sigma^2) + \frac{\|\mathbf{y} - \mathbf{Z}\Lambda(\boldsymbol{\theta})\mathbf{u} - \mathbf{X}\boldsymbol{\beta}\|^2}{\sigma^2} \\
-2\log(h(\mathbf{u}|\mathbf{y}, \boldsymbol{\theta}, \boldsymbol{\beta}, \sigma)) &= (n+q)\log(2\pi\sigma^2) + \frac{\|\mathbf{y} - \mathbf{Z}\Lambda(\boldsymbol{\theta})\mathbf{u} - \mathbf{X}\boldsymbol{\beta}\|^2 + \|\mathbf{u}\|^2}{\sigma^2} \\
&= (n+q)\log(2\pi\sigma^2) + \frac{d(\mathbf{u}|\mathbf{y}, \boldsymbol{\theta}, \boldsymbol{\beta})}{\sigma^2}
\end{aligned} \tag{7}$$

In (7) the *discrepancy* function,

$$d(\mathbf{u}|\mathbf{y}, \boldsymbol{\theta}, \boldsymbol{\beta}) = \|\mathbf{y} - \mathbf{Z}\Lambda(\boldsymbol{\theta})\mathbf{u} - \mathbf{X}\boldsymbol{\beta}\|^2 + \|\mathbf{u}\|^2 \tag{8}$$

has the form of a penalized residual sum of squares in that the first term, $\|\mathbf{y} - \mathbf{Z}\Lambda(\boldsymbol{\theta})\mathbf{u} - \mathbf{X}\boldsymbol{\beta}\|^2$ 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, σ .

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 $\boldsymbol{\beta}$

$$d(\mathbf{u}|\mathbf{y}, \boldsymbol{\theta}, \boldsymbol{\beta}) = \left\| \begin{bmatrix} \mathbf{y} \\ \mathbf{0} \end{bmatrix} - \begin{bmatrix} \mathbf{Z}\Lambda(\boldsymbol{\theta}) & \mathbf{X} \\ \mathbf{I}_q & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{u} \\ \boldsymbol{\beta} \end{bmatrix} \right\|^2. \tag{9}$$

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, $\mathbf{Z}\Lambda(\boldsymbol{\theta})\mathbf{u} + \mathbf{X}\boldsymbol{\beta}$, 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 (9) 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

$$\begin{bmatrix} \Lambda'(\boldsymbol{\theta})\mathbf{Z}'\mathbf{Z}\Lambda(\boldsymbol{\theta}) + \mathbf{I}_q & \Lambda'(\boldsymbol{\theta})\mathbf{Z}\mathbf{X} \\ \mathbf{X}'\mathbf{Z}\Lambda(\boldsymbol{\theta}) & \mathbf{X}'\mathbf{X} \end{bmatrix} \begin{bmatrix} \tilde{\mathbf{u}}(\boldsymbol{\theta}) \\ \tilde{\boldsymbol{\beta}}(\boldsymbol{\theta}) \end{bmatrix} = \begin{bmatrix} \Lambda'(\boldsymbol{\theta})\mathbf{Z}'\mathbf{y} \\ \mathbf{X}'\mathbf{y} \end{bmatrix} \tag{10}$$

An effective way of determining the solution to a sparse, symmetric, positive definite system of equations such as (10) is the sparse Cholesky decomposition (Davis, 2006). If \mathbf{A} is a sparse, symmetric positive definite matrix then the sparse Cholesky decomposition with fill-reducing permutation \mathbf{P} is the lower-triangular matrix \mathbf{L} such that

$$\mathbf{LL}' = \mathbf{PAP}' \quad (11)$$

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

$$\mathbf{P} = \begin{bmatrix} \mathbf{P}_Z & \mathbf{0} \\ \mathbf{0} & \mathbf{P}_X \end{bmatrix} \quad (12)$$

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 (12) so that we can write the Cholesky factorization for the positive definite system (10) as

$$\begin{bmatrix} \mathbf{L}_Z & \mathbf{0} \\ \mathbf{L}_{XZ} & \mathbf{L}_X \end{bmatrix} \begin{bmatrix} \mathbf{L}_Z & \mathbf{0} \\ \mathbf{L}_{XZ} & \mathbf{L}_X \end{bmatrix}' = \begin{bmatrix} \mathbf{P}_Z & \mathbf{0} \\ \mathbf{0} & \mathbf{P}_X \end{bmatrix} \begin{bmatrix} \Lambda'(\boldsymbol{\theta})\mathbf{Z}'\mathbf{Z}\Lambda(\boldsymbol{\theta}) + \mathbf{I}_q & \Lambda'(\boldsymbol{\theta})\mathbf{Z}\mathbf{X} \\ \mathbf{X}'\mathbf{Z}\Lambda(\boldsymbol{\theta}) & \mathbf{X}'\mathbf{X} \end{bmatrix} \begin{bmatrix} \mathbf{P}_Z & \mathbf{0} \\ \mathbf{0} & \mathbf{P}_X \end{bmatrix}' \quad (13)$$

The diagonal elements of \mathbf{L}_Z and \mathbf{L}_X are determined only up to a change in sign. By convention we choose them to be positive.

The discrepancy can now be written in the canonical form

$$d(\mathbf{u}|\mathbf{y}, \boldsymbol{\theta}, \boldsymbol{\beta}) = \tilde{d}(\mathbf{y}, \boldsymbol{\theta}) + \left\| \begin{bmatrix} \mathbf{L}'_Z & \mathbf{L}'_{XZ} \\ \mathbf{0} & \mathbf{L}'_X \end{bmatrix} \begin{bmatrix} \mathbf{P}_Z(\mathbf{u} - \tilde{\mathbf{u}}) \\ \mathbf{P}_X(\boldsymbol{\beta} - \tilde{\boldsymbol{\beta}}) \end{bmatrix} \right\|^2 \quad (14)$$

where

$$\tilde{d}(\mathbf{y}, \boldsymbol{\theta}) = d(\tilde{\mathbf{u}}(\boldsymbol{\theta})|\mathbf{y}, \boldsymbol{\theta}, \tilde{\boldsymbol{\beta}}(\boldsymbol{\theta})) \quad (15)$$

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

3.2 Profiled log-likelihood for linear mixed models

Substituting (14) into (7) provides the unnormalized conditional density $h(\mathbf{u}|\mathbf{y}, \boldsymbol{\theta}, \boldsymbol{\beta}, \sigma)$ on the deviance scale as

$$\begin{aligned} & -2 \log(h(\mathbf{u}|\mathbf{y}, \boldsymbol{\theta}, \boldsymbol{\beta}, \sigma)) \\ &= (n+q) \log(2\pi\sigma^2) + \frac{\tilde{d}(\mathbf{y}, \boldsymbol{\theta}) + \left\| \begin{bmatrix} \mathbf{L}'_{\mathbf{Z}} & \mathbf{L}'_{\mathbf{XZ}} \\ \mathbf{0} & \mathbf{L}'_{\mathbf{X}} \end{bmatrix} \begin{bmatrix} \mathbf{P}_{\mathbf{Z}}(\mathbf{u} - \tilde{\mathbf{u}}) \\ \mathbf{P}_{\mathbf{X}}(\boldsymbol{\beta} - \tilde{\boldsymbol{\beta}}) \end{bmatrix} \right\|^2}{\sigma^2}. \end{aligned} \quad (16)$$

We need to integrate h with respect to \mathbf{u} and minimize the resulting function with respect to the parameter values. We first minimize h with respect to $\boldsymbol{\beta}$ by setting $\boldsymbol{\beta} = \tilde{\boldsymbol{\beta}}(\boldsymbol{\theta})$, from which we can evaluate the *profiled likelihood*, a function of $\boldsymbol{\theta}$ only. On the deviance scale the profiled likelihood is

$$-2 \log(\tilde{L}(\boldsymbol{\theta}|\mathbf{y})) = \log(|\mathbf{L}_{\mathbf{Z}}(\boldsymbol{\theta})|^2) + n \left(1 + \log \left(\frac{2\pi\tilde{d}(\mathbf{y}, \boldsymbol{\theta})}{n} \right) \right) \quad (17)$$

and the maximum likelihood estimate of $\boldsymbol{\theta}$ is

$$\hat{\boldsymbol{\theta}}_L = \arg \min_{\boldsymbol{\theta}} \left(-2 \log(\tilde{L}(\boldsymbol{\theta}|\mathbf{y})) \right). \quad (18)$$

The maximum likelihood estimates of σ^2 and $\boldsymbol{\beta}$ are then evaluated as

$$\hat{\sigma}_L^2 = \frac{\tilde{d}(\hat{\boldsymbol{\theta}}_L, \mathbf{y})}{n} \quad (19)$$

$$\hat{\boldsymbol{\beta}}_L = \tilde{\boldsymbol{\beta}}(\hat{\boldsymbol{\theta}}_L). \quad (20)$$

3.3 The REML criterion

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

$$\mathbf{y} \sim \mathcal{N}(\mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{I}_n), \quad (21)$$

in which we typically estimate σ^2 by

$$\widehat{\sigma_R^2} = \frac{\|\mathbf{y} - \mathbf{X}\widehat{\boldsymbol{\beta}}\|^2}{n - p} \quad (22)$$

even though the maximum likelihood estimate of σ^2 is

$$\widehat{\sigma_L^2} = \frac{\|\mathbf{y} - \mathbf{X}\widehat{\boldsymbol{\beta}}\|^2}{n}. \quad (23)$$

The argument for preferring $\widehat{\sigma_R^2}$ to $\widehat{\sigma_L^2}$ as an estimate of σ^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 σ^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 (22) for a linear regression model. Although not usually derived in this way, the REML criterion can be expressed as

$$c_R(\boldsymbol{\theta}, \boldsymbol{\sigma}|\mathbf{y}) = -2 \log \int_{\mathbb{R}^p} L(\mathbf{u}|\mathbf{y}, \boldsymbol{\theta}, \boldsymbol{\beta}, \boldsymbol{\sigma}) d\boldsymbol{\beta} \quad (24)$$

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

$$\tilde{c}_R(\boldsymbol{\theta}|\mathbf{y}) = \log(|\mathbf{L}\mathbf{z}(\boldsymbol{\theta})|^2) + \log(|\mathbf{L}\mathbf{x}(\boldsymbol{\theta})|^2) + (n - p) \left(1 + \log \left(\frac{2\pi\tilde{d}(\boldsymbol{\theta}|\mathbf{y})}{n - p} \right) \right) \quad (25)$$

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

$$\widehat{\boldsymbol{\theta}}_R = \arg \min_{\boldsymbol{\theta}} \tilde{c}_R(\boldsymbol{\theta}, \mathbf{y}). \quad (26)$$

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

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 (24), does not depend on

β . However, it is customary (and not unreasonable) to use $\hat{\beta}_R = \tilde{\beta}(\hat{\theta}_r)$ as the REML estimate of β .

Note that the profiled REML criterion can be evaluated from a sparse Cholesky decomposition like that in (13) but without the requirement that the permutation can be applied to the columns of $\mathbf{Z}\Lambda(\theta)$ separately from the columns of \mathbf{X} . That is, we can use a general fill-reducing permutation rather than the specific form (12) 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 Implementation details for linear mixed models

We determine the maximum likelihood (ML) or REML estimates of the parameters θ , β and σ in a linear mixed model by first minimizing the profiled log-likelihood or the profiled REML criterion with respect to θ . Thus the first implementation question is evaluating the profiled log-likelihood or the profiled REML criterion.

For many models, in particular models with scalar random effects, which are described later, the matrix $\Lambda(\theta)$ 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

$$\mathbf{A} = \begin{bmatrix} \mathbf{Z}'\mathbf{Z} & \mathbf{Z}'\mathbf{X} \\ \mathbf{X}'\mathbf{Z} & \mathbf{X}'\mathbf{X} \end{bmatrix} \quad \text{and} \quad \mathbf{c} = \begin{bmatrix} \mathbf{Z}'\mathbf{y} \\ \mathbf{X}'\mathbf{y} \end{bmatrix} \quad (27)$$

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

$$\mathbf{L}(\theta)\mathbf{L}(\theta)' = \mathbf{P} \left(\begin{bmatrix} \Lambda(\theta) & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_p \end{bmatrix} \mathbf{A} \begin{bmatrix} \Lambda(\theta) & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_p \end{bmatrix} + \begin{bmatrix} \mathbf{I}_q & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \right) \mathbf{P}' \quad (28)$$

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

$$\mathbf{L}\mathbf{L}'\mathbf{P} \begin{bmatrix} \tilde{\mathbf{u}}(\theta) \\ \tilde{\beta}(\theta) \end{bmatrix} = \mathbf{P} \begin{bmatrix} \Lambda(\theta) & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_p \end{bmatrix} \mathbf{c} \quad (29)$$

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

$$\tilde{d}_R(\theta|\mathbf{y}) = \log(|\mathbf{L}(\theta)|^2) + (n-p) \left(1 + \log \left(\frac{2\pi\tilde{d}(\mathbf{y}|\theta)}{n-p} \right) \right). \quad (30)$$

References

- 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, $\boldsymbol{\Lambda}(\boldsymbol{\theta})$ of the relative covariance matrix of \mathbf{B} (dimension m).
- σ 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

L Left Cholesky factor of a positive-definite symmetric matrix. L_Z is $q \times q$; L_X is $p \times p$.

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

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

X Model matrix for the fixed-effects parameters, β . (Size $(ns) \times p$.)

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

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

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.