

See discussions, stats, and author profiles for this publication at: <https://www.researchgate.net/publication/2753537>

# Computational Methods for Multilevel Modelling

Article · March 1998

Source: CiteSeer

---

CITATIONS

56

---

READS

1,405

2 authors, including:



[Douglas M. Bates](#)

University of Wisconsin–Madison

137 PUBLICATIONS 161,979 CITATIONS

SEE PROFILE

# Computational Methods for Multilevel Modelling

Douglas M. Bates

Department of Statistics  
University of Wisconsin – Madison \*

José C. Pinheiro

Bell Laboratories  
Lucent Technologies

## Abstract

A multilevel mixed-effects model has random effects at each of several nested levels of grouping of the observed responses. We may use these, for example, when modelling observations taken over time on students who are grouped into classes that are grouped into schools that are grouped into districts. If each of the distributions of the random effects is Gaussian and if the disturbance term at the lowest level of grouping is also Gaussian it is straightforward to define a likelihood for the fixed effects and the parameters defining the random effects distribution. We show that by expressing the random effects distribution in terms of relative precision factors and using matrix decompositions, this likelihood can be profiled and can be compactly expressed. The same decompositions provide rapid evaluation of the profiled log-restricted-likelihood for REML estimation.

The conditional distribution of the random effects given the data can be derived from the decomposed matrices. From this a compact and rapidly evaluated expression for the EM iterations can be derived. Reasonable starting estimates for the relative precision factors can be derived from the design alone. These starting estimates, refined by a moderate number of EM iterations, provide excellent starting values for a Newton-Raphson or quasi-Newton optimization of the log-likelihood or the log-restricted-likelihood. The methods we describe extend easily to models with non-spherical distributions for the within-group errors and to nonlinear multilevel models.

*Key words and phrases:* mixed-effects models, EM algorithm, maximum likelihood, restricted maximum likelihood

## 1 Introduction

We consider computational methods for Gaussian multilevel mixed-effects models as described, for example, in Longford (1993) or Goldstein (1995). These models are used with data where the individual observations are grouped at one or more hierarchical levels. For example, we may wish to model observations on students who are grouped into

---

\*This research was supported by the National Science Foundation through grant DMS-9704349.

classes that are grouped into schools that are grouped into school districts. If, at each level in the grouping hierarchy, the experimental units that are observed constitute a sample from the population about which we wish to make inferences, we model the effect of each unit at each level as a random effect. A model including random effects for the individual units as well as overall fixed effects is called a mixed-effects model.

Our computational methods are suitable for maximum likelihood or maximum restricted likelihood estimation of the parameters in linear mixed-effects models with Gaussian distributions for the random effects at each level. These techniques can also form the basis for parameter estimation methods for general linear mixed-effects models or for nonlinear mixed-effects models.

## 1.1 Description of the model

For a single level of grouping, the linear mixed-effects model described by Laird and Ware (1982) expresses the  $n_i$ -dimensional response vector  $\mathbf{y}_i$  for the  $i$ th unit as

$$\begin{aligned}\mathbf{y}_i &= \mathbf{X}_i\boldsymbol{\beta} + \mathbf{Z}_i\mathbf{b}_i + \boldsymbol{\epsilon}_i, \quad i = 1, \dots, M \\ \mathbf{b}_i &\sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}), \quad \boldsymbol{\epsilon}_i \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})\end{aligned}\tag{1}$$

where  $\boldsymbol{\beta}$  is the  $p$ -dimensional vector of *fixed effects*,  $\mathbf{b}_i$  is the  $q$ -dimensional vector of *random effects*,  $\mathbf{X}_i$  (of size  $n_i \times p$ ) and  $\mathbf{Z}_i$  (of size  $n_i \times q$ ) are known fixed-effects and random-effects regressor matrices, and  $\boldsymbol{\epsilon}_i$  is the  $n_i$ -dimensional *within-group error* vector with a spherical Gaussian distribution. The assumption  $\text{Var}(\boldsymbol{\epsilon}_i) = \sigma^2 \mathbf{I}$  can be relaxed as shown in §5.1.

The random effects for the  $i$ th unit,  $\mathbf{b}_i$ , are also assumed to have a Gaussian distribution but with a general positive-definite variance-covariance matrix  $\boldsymbol{\Sigma}$ .

In some of the multilevel modelling literature, notably in Goldstein (1995), the model (1) is called a “two-level model” because there are two levels of random variation:  $\mathbf{b}_i$  and  $\boldsymbol{\epsilon}_i$ . In other references this model would be described as having one level of random effects in the model. We will adopt the latter convention and count the “levels” in a multilevel model as the number of levels of nested random effects.

In a multilevel model with two levels of random effects the  $n_{ij}$ -dimensional vector of responses for the  $j$ th level-2 unit nested within the  $i$ th level-1 unit is written

$$\begin{aligned}\mathbf{y}_{ij} &= \mathbf{X}_{ij}\boldsymbol{\beta} + \mathbf{Z}_{i,j}\mathbf{b}_i + \mathbf{Z}_{ij}\mathbf{b}_{ij} + \boldsymbol{\epsilon}_{ij} \quad i = 1, \dots, M \quad j = 1, \dots, M_i \\ \mathbf{b}_i &\sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}_1), \quad \mathbf{b}_{ij} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}_2), \quad \boldsymbol{\epsilon}_{ij} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})\end{aligned}\tag{2}$$

The regressor matrices  $\mathbf{X}_{ij}$  (of size  $n_{ij} \times p$ ),  $\mathbf{Z}_{i,j}$  (of size  $n_{ij} \times q_1$ ) and  $\mathbf{Z}_{ij}$  (of size  $n_{ij} \times q_2$ ) correspond to the fixed effects  $\boldsymbol{\beta}$  and the first- and second-level random effects  $\mathbf{b}_i$  and  $\mathbf{b}_{ij}$ .

Extensions to an arbitrary number of levels of random effects follow the same general pattern. For example, with three levels of random effects the response for the  $k$ th level-3 unit within the  $j$ th level-2 unit within the  $i$ th level-1 unit will be written

$$\begin{aligned}\mathbf{y}_{ijk} &= \mathbf{X}_{ijk}\boldsymbol{\beta} + \mathbf{Z}_{i,j,k}\mathbf{b}_i + \mathbf{Z}_{ij,k}\mathbf{b}_{ij} + \mathbf{Z}_{ijk}\mathbf{b}_{ijk} + \boldsymbol{\epsilon}_{ijk} \\ i &= 1, \dots, M \quad j = 1, \dots, M_i \quad k = 1, \dots, M_{ij} \\ \mathbf{b}_i &\sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}_1), \quad \mathbf{b}_{ij} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}_2), \quad \mathbf{b}_{ijk} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}_3), \quad \boldsymbol{\epsilon}_{ijk} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})\end{aligned}$$

Note that the distinction between, say, the  $k$ th horizontal section of the regressor matrix for the level-2 random effect  $b_{ij}$ , written  $Z_{ij,k}$ , and the  $j$ th horizontal section of the regressor matrix for the level-1 random effect  $b_i$ , written  $Z_{i,jk}$ , is the position of the comma in the subscripts.

## 1.2 The variance-covariance of the random effects

In a model with  $Q$  nested levels of random effects, the variance-covariance matrices  $\Sigma_q$ ,  $q = 1, \dots, Q$  of the random effects must be symmetric and at least positive semi-definite. We consider only positive definite  $\Sigma_q$  because any indefiniteness can be removed by re-expressing the model in terms of a random-effects vector of lower dimension.

Within the set of positive-definite  $\Sigma_q$  we may impose further constraints such as requiring that  $\Sigma_q$  be diagonal or that it be a multiple of the identity or that it have a particular structure such as compound symmetry.

Whatever structure we assume for  $\Sigma_q$  we will always express  $\Sigma_q$  as a function of an unconstrained parameter vector  $\theta_q$  as described in Pinheiro and Bates (1996). By doing this we can ensure that the optimization of the likelihood can be expressed as an unconstrained optimization.

We extend the techniques of Pinheiro and Bates (1996) by expressing the conversion from  $\theta_q$  to  $\Sigma_q$  is two stages. First, we write

$$\Sigma_q = \sigma^2 D_q$$

where  $\sigma^2$  is the variance of the components of  $\epsilon_i$ . The matrix  $D_q$  is thus a relative or scaled variance-covariance matrix for the random effects. This re-expression of  $\Sigma_q$  does not change the form of the model but, as shown in §2, it does provide a more convenient expression of the likelihood or restricted likelihood for the model.

The next stage is to use a “square-root” factor of the relative precision matrix  $D_q^{-1}$ . This *relative precision factor*,  $\Delta_q$ , is any matrix such that

$$D_q^{-1} = \Delta_q' \Delta_q$$

Such a  $\Delta_q$  always exists but does not have to be unique. We could, for example, use the Cholesky factor (Thisted, 1988, §3.3) of  $D_q^{-1}$  or the transpose of the inverse of the Cholesky factor of  $D_q$ .

The parameter vector  $\theta_q$  defines  $\Delta_q$  through some unconstrained parameterization from which  $D_q$  and  $\Sigma_q$  are subsequently defined. Especially in deriving the EM iterations (§3.4) we can work with  $\Delta_q$  directly so the parameterization used is not an issue. When we do need to work in terms of  $\theta_q$  one of the parameterizations considered in Pinheiro and Bates (1996) could be used to define  $\Delta_q$  from  $\theta_q$ .

## 1.3 An example: variability in IC manufacturing

To illustrate the multilevel model described in this section and the methodology presented in the following sections, we present an example from integrated circuit (IC) manufacturing.

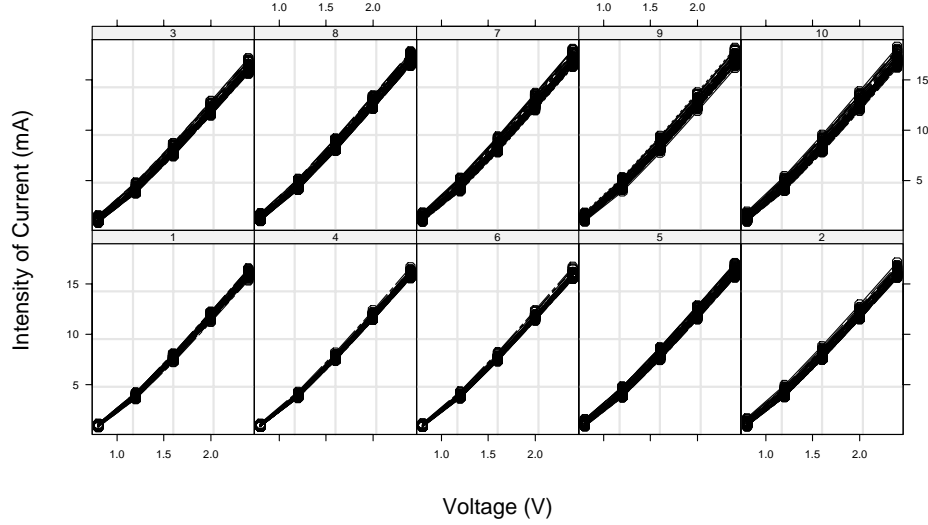


Figure 1: Current versus voltage curves for each of 67 sites within 10 wafers. Each panel presents data for the wafer whose number is given above the panel. The panels are ordered starting from the lower left by increasing maximum current intensity. There are 67 lines, one for each site, on each panel.

In an experiment conducted at the Microelectronics Division of Lucent Technologies to study different sources of variability in the manufacturing of analog MOS circuits, the intensity of current (in milli-Amperes) at 0.8, 1.2, 1.6, 2.0, and 2.4 Volts was measured on  $80\mu\text{m} \times 0.6\mu\text{m}$  n-channel devices. Measurements were made on 10 wafers, each subdivided into 67 sites containing one device. The data are presented on Figure 1, where each panel represents a different wafer and each curve on a panel represents a different site.

Two levels of nesting are present in these data: *wafer* and *site within wafer*. The main objective of the experiment was to construct an empirical model for simulating the behavior of similar circuits.

In Figure 1 it appears that current could be modelled as a quadratic function of voltage. Preliminary analyses indicate that random effects are needed to account for the wafer-to-wafer variability of the intercept and the linear terms as well as for the site-to-site variability of the intercept. The corresponding multilevel model for the intensities of current in the  $j$ th site within the  $i$ th wafer is expressed, for  $i = 1, \dots, 10$  and

$j = 1, \dots, 67$ , as

$$\underbrace{\begin{bmatrix} y_{ij_1} \\ y_{ij_2} \\ y_{ij_3} \\ y_{ij_4} \\ y_{ij_5} \end{bmatrix}}_{\mathbf{y}_{ij}} = \underbrace{\begin{bmatrix} 1 & 0.8 & 0.8^2 \\ 1 & 1.2 & 1.2^2 \\ 1 & 1.6 & 1.6^2 \\ 1 & 2.0 & 2.0^2 \\ 1 & 2.4 & 2.4^2 \end{bmatrix}}_{\mathbf{X}_{ij}} \underbrace{\begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \end{bmatrix}}_{\boldsymbol{\beta}} + \underbrace{\begin{bmatrix} 1 & 0.8 \\ 1 & 1.2 \\ 1 & 1.6 \\ 1 & 2.0 \\ 1 & 2.4 \end{bmatrix}}_{\mathbf{Z}_{i,j}} \underbrace{\begin{bmatrix} b_{i1} \\ b_{i2} \end{bmatrix}}_{\mathbf{b}_i} + \underbrace{\begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \\ 1 \end{bmatrix}}_{\mathbf{Z}_{ij}} \underbrace{\begin{bmatrix} b_{ij_1} \end{bmatrix}}_{\mathbf{b}_{ij}} + \underbrace{\begin{bmatrix} \epsilon_{ij_1} \\ \epsilon_{ij_2} \\ \epsilon_{ij_3} \\ \epsilon_{ij_4} \\ \epsilon_{ij_5} \end{bmatrix}}_{\boldsymbol{\epsilon}_{ij}}$$

$$\begin{bmatrix} b_{i1} \\ b_{i2} \end{bmatrix} \sim \mathcal{N} \left( \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} \Sigma_{111} & \Sigma_{112} \\ \Sigma_{112} & \Sigma_{122} \end{bmatrix} \right), \quad b_{ij_1} \sim \mathcal{N}(0, \Sigma_2), \quad \epsilon_{ij_k} \sim \mathcal{N}(0, \sigma^2)$$

In this example  $Q = 2$ ,  $M = 10$ ,  $M_i = 67$ ,  $n_{ij} = 5$ ,  $q_1 = 2$ ,  $q_2 = 1$ ,  $\boldsymbol{\theta}_1$  is of dimension 3, and  $\boldsymbol{\theta}_2$  is of dimension 1.

## 2 Expressing the likelihood

### 2.1 Single level of random effects

The likelihood function for model (1) can be written

$$\begin{aligned} L(\boldsymbol{\beta}, \boldsymbol{\theta}, \sigma^2 | \mathbf{y}) &= \prod_{i=1}^M p(\mathbf{y}_i | \boldsymbol{\beta}, \boldsymbol{\theta}, \sigma^2) \\ &= \prod_{i=1}^M \int p(\mathbf{y}_i | \mathbf{b}_i, \boldsymbol{\beta}, \sigma^2) p(\mathbf{b}_i | \boldsymbol{\theta}, \sigma^2) d\mathbf{b}_i \\ &= \prod_{i=1}^M \frac{1}{\sqrt{(2\pi\sigma^2)^{n_i} |\mathbf{D}|}} \times \\ &\quad \int \frac{\exp \left[ \frac{-1}{2\sigma^2} (\|\mathbf{y}_i - \mathbf{X}_i \boldsymbol{\beta} - \mathbf{Z}_i \mathbf{b}_i\|^2 + \mathbf{b}_i' \mathbf{D}^{-1} \mathbf{b}_i) \right]}{(2\pi\sigma^2)^{q/2}} d\mathbf{b}_i \quad (3) \end{aligned}$$

The expression  $\|\mathbf{y}_i - \mathbf{X}_i \boldsymbol{\beta} - \mathbf{Z}_i \mathbf{b}_i\|^2 + \mathbf{b}_i' \mathbf{D}^{-1} \mathbf{b}_i$  in the exponent within the integral has the form of a penalized residual sum-of-squares. The term  $\|\mathbf{y}_i - \mathbf{X}_i \boldsymbol{\beta} - \mathbf{Z}_i \mathbf{b}_i\|^2$  is exactly the residual sum of squares for the  $i$ th unit, and the additional term,  $\mathbf{b}_i' \mathbf{D}^{-1} \mathbf{b}_i$ , can be viewed as a “penalty” that inhibits the size of the random-effects vector  $\mathbf{b}_i$ . Using a relative precision factor  $\boldsymbol{\Delta}$  we can write the penalty term in the form of a residual sum-of-squares as

$$\mathbf{b}_i' \mathbf{D}^{-1} \mathbf{b}_i = \|\boldsymbol{\Delta} \mathbf{b}_i\|^2 = \|\mathbf{0} - \mathbf{0} \boldsymbol{\beta} - \boldsymbol{\Delta} \mathbf{b}_i\|^2$$

This re-expression of the penalty is sometimes called a “pseudo-data” representation because it is equivalent to augmenting the observations and the regressor matrices with  $q$  new rows that look like additional observations. Writing

$$\tilde{\mathbf{y}}_i = \begin{bmatrix} \mathbf{y}_i \\ \mathbf{0} \end{bmatrix} \quad \tilde{\mathbf{X}}_i = \begin{bmatrix} \mathbf{X}_i \\ \mathbf{0} \end{bmatrix} \quad \tilde{\mathbf{Z}}_i = \begin{bmatrix} \mathbf{Z}_i \\ \boldsymbol{\Delta} \end{bmatrix}$$

the expression in the exponent becomes

$$\begin{aligned}\|\mathbf{y}_i - \mathbf{X}_i\boldsymbol{\beta} - \mathbf{Z}_i\mathbf{b}_i\|^2 + \mathbf{b}_i' \mathbf{D}^{-1} \mathbf{b}_i &= \|\mathbf{y}_i - \mathbf{X}_i\boldsymbol{\beta} - \mathbf{Z}_i\mathbf{b}_i\|^2 + \|\mathbf{0} - \mathbf{0}\boldsymbol{\beta} - \boldsymbol{\Delta}\mathbf{b}_i\|^2 \\ &= \|\tilde{\mathbf{y}}_i - \tilde{\mathbf{X}}_i\boldsymbol{\beta} - \tilde{\mathbf{Z}}_i\mathbf{b}_i\|^2\end{aligned}$$

Having made the exponent expression look like a residual sum-of-squares we can employ standard numerical techniques for least squares problems. If we form the orthogonal-triangular decomposition (Thisted, 1988, §3.1)

$$\tilde{\mathbf{Z}}_i = \mathbf{Q}_{(i)} \begin{bmatrix} \mathbf{R}_{11(i)} \\ \mathbf{0} \end{bmatrix}$$

where  $\mathbf{Q}_{(i)}$  is a  $(n_i + q) \times (n_i + q)$  orthogonal matrix and  $\mathbf{R}_{11(i)}$  is an upper-triangular  $q \times q$  matrix, then the properties of orthogonal matrices ensure that

$$\begin{aligned}\|\tilde{\mathbf{y}}_i - \tilde{\mathbf{X}}_i\boldsymbol{\beta} - \tilde{\mathbf{Z}}_i\mathbf{b}_i\|^2 &= \|\mathbf{Q}'_{(i)} (\tilde{\mathbf{y}}_i - \tilde{\mathbf{X}}_i\boldsymbol{\beta} - \tilde{\mathbf{Z}}_i\mathbf{b}_i)\|^2 \\ &= \|\mathbf{c}_{1(i)} - \mathbf{R}_{10(i)}\boldsymbol{\beta} - \mathbf{R}_{11(i)}\mathbf{b}_i\|^2 + \|\mathbf{c}_{0(i)} - \mathbf{R}_{00(i)}\boldsymbol{\beta}\|^2\end{aligned}$$

where the  $q \times p$  matrix  $\mathbf{R}_{10(i)}$ , the  $n_i \times p$  matrix  $\mathbf{R}_{00(i)}$ , the  $q$ -vector  $\mathbf{c}_{1(i)}$  and  $n_i$ -vector  $\mathbf{c}_{0(i)}$  are defined by

$$\begin{bmatrix} \mathbf{R}_{10(i)} \\ \mathbf{R}_{00(i)} \end{bmatrix} = \mathbf{Q}'_{(i)} \tilde{\mathbf{X}}_i \quad \text{and} \quad \begin{bmatrix} \mathbf{c}_{1(i)} \\ \mathbf{c}_{0(i)} \end{bmatrix} = \mathbf{Q}'_{(i)} \tilde{\mathbf{y}}_i$$

Furthermore, if  $\boldsymbol{\Sigma}$  is positive definite, as we require, then  $\boldsymbol{\Delta}$  is non-singular and hence  $\mathbf{R}_{11(i)}$  is also non-singular.

Another way of thinking of this decomposition is as the orthogonal-triangular (QR) decomposition of an augmented matrix

$$\begin{bmatrix} \mathbf{Z}_i & \mathbf{X}_i & \mathbf{y}_i \\ \boldsymbol{\Delta} & \mathbf{0} & \mathbf{0} \end{bmatrix} = \mathbf{Q}_{(i)} \begin{bmatrix} \mathbf{R}_{11(i)} & \mathbf{R}_{10(i)} & \mathbf{c}_{1(i)} \\ \mathbf{0} & \mathbf{R}_{00(i)} & \mathbf{c}_{0(i)} \end{bmatrix} \quad (4)$$

where the reduction to triangular form is halted after the first  $q$  columns. (The peculiar numbering scheme for the submatrices and subvectors is designed to allow easy extension to more than one level of random effects as seen in §2.4 and §2.5.)

The calculation of the decomposition in (4) is straightforward, efficient, and numerically stable. Standard software such as Linpack (Dongarra, Bunch, Moler and Stewart, 1979) or LAPACK (Anderson, Bai, Bischoff, Demmel, Dongarra, DuCroz, Greenbaum, Hammarling, McKenney, Ostrouchov and Sorensen, 1994) can be used.

Returning to the integral in (3) we can now remove a constant factor and reduce it to

$$\begin{aligned}& \int \frac{\exp \left[ \frac{-1}{2\sigma^2} (\|\mathbf{y}_i - \mathbf{X}_i\boldsymbol{\beta} - \mathbf{Z}_i\mathbf{b}_i\|^2 + \mathbf{b}_i' \mathbf{D}^{-1} \mathbf{b}_i) \right]}{(2\pi\sigma^2)^{q/2}} d\mathbf{b}_i \\ &= \exp \left[ \frac{\|\mathbf{c}_{0(i)} - \mathbf{R}_{00(i)}\boldsymbol{\beta}\|^2}{-2\sigma^2} \right] \int \frac{\exp \left[ \frac{-1}{2\sigma^2} (\|\mathbf{c}_{1(i)} - \mathbf{R}_{10(i)}\boldsymbol{\beta} - \mathbf{R}_{11(i)}\mathbf{b}_i\|^2) \right]}{(2\pi\sigma^2)^{q/2}} d\mathbf{b}_i\end{aligned} \quad (5)$$

Because  $\mathbf{R}_{11(i)}$  is non-singular we can perform a change of variable to  $\phi_i = (\mathbf{c}_{1(i)} - \mathbf{R}_{10(i)}\boldsymbol{\beta} - \mathbf{R}_{11(i)}\mathbf{b}_i)/\sigma$  with the differential  $d\phi_i = \sigma^{-q} \text{abs}|\mathbf{R}_{11(i)}| d\mathbf{b}_i$  and write the integral as

$$\begin{aligned} \int \frac{\exp\left[\frac{-1}{2\sigma^2} (\|\mathbf{c}_{1(i)} - \mathbf{R}_{10(i)}\boldsymbol{\beta} - \mathbf{R}_{11(i)}\mathbf{b}_i\|^2)\right]}{(2\pi\sigma^2)^{q/2}} d\mathbf{b}_i \\ = \frac{1}{\text{abs}|\mathbf{R}_{11(i)}|} \int \frac{\exp(-\|\phi_i\|^2/2)}{(2\pi)^{q/2}} d\phi_i \\ = \text{abs}|\mathbf{R}_{11(i)}|^{-1} \end{aligned} \quad (6)$$

Since, by construction,  $\mathbf{R}_{11(i)}$  is upper-triangular, its determinant is simply the product of its diagonal elements.

Substituting (6) into (5) into (3) provides the likelihood as

$$L(\boldsymbol{\beta}, \mathbf{D}, \sigma^2 | \mathbf{y}) = \prod_{i=1}^M \frac{\exp\left(\frac{\|\mathbf{c}_{0(i)} - \mathbf{R}_{00(i)}\boldsymbol{\beta}\|^2}{-2\sigma^2}\right)}{\sqrt{(2\pi\sigma^2)^{n_i} |\mathbf{D}|}} \text{abs}|\mathbf{R}_{11(i)}|^{-1}$$

which is now in the form of a regression model for the fixed effects  $\boldsymbol{\beta}$ . Forming another orthogonal-triangular decomposition

$$\begin{bmatrix} \mathbf{R}_{00(1)} & \mathbf{c}_{0(1)} \\ \vdots & \vdots \\ \mathbf{R}_{00(M)} & \mathbf{c}_{0(M)} \end{bmatrix} = \mathbf{Q}_0 \begin{bmatrix} \mathbf{R}_{00} & \mathbf{c}_0 \\ \mathbf{0} & \mathbf{c}_{-1} \end{bmatrix} \quad (7)$$

and noting that  $1/\sqrt{|\mathbf{D}|} = \text{abs}|\boldsymbol{\Delta}|$  reduces this to

$$\begin{aligned} L(\boldsymbol{\beta}, \boldsymbol{\theta}, \sigma^2 | \mathbf{y}) \\ = (2\pi\sigma^2)^{-N/2} \exp\left(\frac{\|\mathbf{c}_{-1}\|^2 + \|\mathbf{c}_0 - \mathbf{R}_{00}\boldsymbol{\beta}\|^2}{-2\sigma^2}\right) \prod_{i=1}^M \text{abs}\left(\frac{|\boldsymbol{\Delta}|}{|\mathbf{R}_{11(i)}|}\right) \end{aligned} \quad (8)$$

where  $N = \sum_{i=1}^M n_i$  is the total number of observations.

## 2.2 The profiled log-likelihood for $\boldsymbol{\theta}$

With the likelihood expressed as (8) we can derive explicit expressions for the optimal values of  $\boldsymbol{\beta}$  and  $\sigma^2$  conditional on a value of  $\boldsymbol{\theta}$ . We assume  $\mathbf{R}_{00}$  is non-singular, in which case the optimal value of  $\boldsymbol{\beta}$  satisfies

$$\mathbf{R}_{00}\hat{\boldsymbol{\beta}}(\boldsymbol{\theta}) = \mathbf{c}_0 \quad (9)$$

Substituting this value in (8) and taking the logarithm provides the profiled log-likelihood

$$\begin{aligned} \ell(\boldsymbol{\theta}, \sigma^2 | \mathbf{y}) &= \log L(\hat{\boldsymbol{\beta}}(\boldsymbol{\theta}), \boldsymbol{\theta}, \sigma^2 | \mathbf{y}) \\ &= -\frac{N \log(2\pi\sigma^2)}{2} + \frac{\|\mathbf{c}_{-1}\|^2}{-2\sigma^2} + \sum_{i=1}^M \log \text{abs}\left(\frac{|\boldsymbol{\Delta}|}{|\mathbf{R}_{11(i)}|}\right) \end{aligned}$$



which is maximized with respect to  $\sigma^2$  by  $\widehat{\sigma^2}(\boldsymbol{\theta}) = \|\mathbf{c}_{-1}\|^2/N$

We can now write the profiled log-likelihood as a function of  $\boldsymbol{\theta}$  alone as

$$\begin{aligned}\ell(\boldsymbol{\theta}|\mathbf{y}) &= \log L(\widehat{\boldsymbol{\beta}}(\boldsymbol{\theta}), \boldsymbol{\theta}, \widehat{\sigma^2}(\boldsymbol{\theta})|\mathbf{y}) \\ &= \text{const} - N \log \|\mathbf{c}_{-1}\| + \sum_{i=1}^M \log \text{abs} \left( \frac{|\boldsymbol{\Delta}|}{|\mathbf{R}_{11(i)}|} \right)\end{aligned}\quad (10)$$

### 2.3 Restricted log-likelihood as a function of $\boldsymbol{\theta}$ alone

The *restricted likelihood* (Harville, 1976) is often preferred to the likelihood when defining an estimator for  $\boldsymbol{\theta}$ . One way of writing the restricted likelihood is

$$L_R(\boldsymbol{\theta}, \sigma^2|\mathbf{y}) = \int L(\boldsymbol{\beta}, \boldsymbol{\theta}, \sigma^2|\mathbf{y}) d\boldsymbol{\beta}$$

which reduces to

$$L_R(\boldsymbol{\theta}, \sigma^2|\mathbf{y}) = (2\pi\sigma^2)^{(N-p)/2} \exp\left(\frac{\|\mathbf{c}_{-1}\|^2}{-2\sigma^2}\right) \text{abs} |\mathbf{R}_{00}|^{-1} \prod_{i=1}^M \text{abs} \left( \frac{|\boldsymbol{\Delta}|}{|\mathbf{R}_{11(i)}|} \right)$$

using (8) and the same change-of-variable technique used to obtain (6).

Converting to the log-restricted-likelihood

$$\begin{aligned}\ell_R(\boldsymbol{\theta}, \sigma^2|\mathbf{y}) &= \\ &= -\frac{N-p}{2} \log(2\pi\sigma^2) - \frac{\|\mathbf{c}_{-1}\|^2}{2\sigma^2} - \log \text{abs} |\mathbf{R}_{00}| + \sum_{i=1}^M \log \text{abs} \left( \frac{|\boldsymbol{\Delta}|}{|\mathbf{R}_{11(i)}|} \right)\end{aligned}$$

provides the conditional estimate  $\widehat{\sigma_R^2}(\boldsymbol{\theta}) = \|\mathbf{c}_{-1}\|^2/(N-p)$  for  $\sigma^2$  from which we obtain the profiled log-restricted-likelihood

$$\begin{aligned}\ell_R(\boldsymbol{\theta}|\mathbf{y}) &= \ell_R(\boldsymbol{\theta}, \widehat{\sigma_R^2}(\boldsymbol{\theta})) \\ &= \text{const} - (N-p) \log \|\mathbf{c}_{-1}\| - \log \text{abs} |\mathbf{R}_{00}| + \sum_{i=1}^M \log \text{abs} \left( \frac{|\boldsymbol{\Delta}|}{|\mathbf{R}_{11(i)}|} \right)\end{aligned}$$

### 2.4 Two levels of random effects

The likelihood for a model with two levels of random effects is defined as in (3) but integrating over both levels of random effects

$$\begin{aligned}L(\boldsymbol{\beta}, \boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \sigma^2|\mathbf{y}) &= \\ &= \prod_{i=1}^M \int \prod_{j=1}^{M_i} \left[ \int p(\mathbf{y}_{ij}|\mathbf{b}_{ij}, \mathbf{b}_i, \boldsymbol{\beta}, \sigma^2) p(\mathbf{b}_{ij}|\boldsymbol{\theta}_2, \sigma^2) d\mathbf{b}_{ij} \right] p(\mathbf{b}_i|\boldsymbol{\theta}_1, \sigma^2) d\mathbf{b}_i\end{aligned}\quad (11)$$

As with the single level of random effects, we can simplify the integrals in (11) if we augment the  $\mathbf{Z}_{ij}$  matrices with  $\Delta_2$  and form orthogonal-triangular decompositions of these augmented arrays. This allows us to evaluate the inner integrals. To evaluate the outer integrals we iterate this process.

That is, we first form and decompose the arrays

$$\begin{bmatrix} \mathbf{Z}_{ij} & \mathbf{Z}_{i,j} & \mathbf{X}_{ij} & \mathbf{y}_{ij} \\ \Delta_2 & \mathbf{0} & \mathbf{0} & \mathbf{0} \end{bmatrix} = \mathbf{Q}_{(ij)} \begin{bmatrix} \mathbf{R}_{22(ij)} & \mathbf{R}_{21(ij)} & \mathbf{R}_{20(ij)} & \mathbf{c}_{2(ij)} \\ \mathbf{0} & \mathbf{R}_{11(ij)} & \mathbf{R}_{10(ij)} & \mathbf{c}_{1(ij)} \end{bmatrix} \quad (12)$$

$$i = 1, \dots, M \quad j = 1, \dots, M_i$$

The matrix  $\mathbf{R}_{22(ij)}$  will be an upper-triangular matrix of dimension  $q_2 \times q_2$ . The other arrays in the first row of the decomposition in (12) are used only if the conditional estimates of  $\beta$  or the Best Linear Unbiased Predictors (BLUPs) (see §3.5) for  $\mathbf{b}_{ij}$  and  $\mathbf{b}_i$  are required. The arrays in the second row of the decomposition:  $\mathbf{R}_{11(ij)}$ ,  $\mathbf{R}_{10ij}$ , and  $\mathbf{c}_{1(ij)}$  each have  $n_{ij}$  rows.

To evaluate the outer integral in (11) we again form and decompose an augmented array

$$\begin{bmatrix} \mathbf{R}_{11(i1)} & \mathbf{R}_{10(i1)} & \mathbf{c}_{1(i1)} \\ \vdots & \vdots & \vdots \\ \mathbf{R}_{11(iM_i)} & \mathbf{R}_{10(iM_i)} & \mathbf{c}_{1(iM_i)} \\ \Delta_1 & \mathbf{0} & \mathbf{0} \end{bmatrix} = \mathbf{Q}_{(i)} \begin{bmatrix} \mathbf{R}_{11(i)} & \mathbf{R}_{10(i)} & \mathbf{c}_{1(i)} \\ \mathbf{0} & \mathbf{R}_{00(i)} & \mathbf{c}_{0(i)} \end{bmatrix} \quad i = 1, \dots, M \quad (13)$$

The final decomposition to produce  $\mathbf{R}_{00}$ ,  $\mathbf{c}_0$  and  $\mathbf{c}_{-1}$  is the same as that in (7).

Using the matrices and vectors produced in (12), (13), and (7) and following the same steps as for the single level of nesting we can express the profiled log-likelihood for  $\theta_1$  and  $\theta_2$  as

$$\begin{aligned} \ell(\theta_1, \theta_2 | \mathbf{y}) &= \log L(\hat{\beta}(\theta_1, \theta_2), \theta_1, \theta_2, \widehat{\sigma^2}(\theta_1, \theta_2) | \mathbf{y}) \\ &= \text{const} - N \log \|\mathbf{c}_{-1}\| + \sum_{i=1}^M \log \text{abs} \left( \frac{|\Delta_1|}{|\mathbf{R}_{11(i)}|} \right) \\ &\quad + \sum_{i=1}^M \sum_{j=1}^{M_i} \log \text{abs} \left( \frac{|\Delta_2|}{|\mathbf{R}_{22(ij)}|} \right) \end{aligned}$$

Similarly, the profiled log-restricted-likelihood is

$$\begin{aligned} \ell_R(\theta_1, \theta_2 | \mathbf{y}) &= \log L_R(\hat{\beta}_R(\theta_1, \theta_2), \theta_1, \theta_2, \widehat{\sigma^2}_R(\theta_1, \theta_2) | \mathbf{y}) \\ &= \text{const} - (N - p) \log \|\mathbf{c}_{-1}\| - \log \text{abs} |\mathbf{R}_{00}| \\ &\quad + \sum_{i=1}^M \log \text{abs} \left( \frac{|\Delta_1|}{|\mathbf{R}_{11(i)}|} \right) + \sum_{i=1}^M \sum_{j=1}^{M_i} \log \text{abs} \left( \frac{|\Delta_2|}{|\mathbf{R}_{22(ij)}|} \right) \end{aligned}$$

To illustrate the decomposition in (12) and the first part of the calculation of the likelihood, we consider the data from the first site within the first wafer in the IC manufacturing example from §1.3. As described in §2.6 we use 0.84 as an initial value for

$\Delta_2$ . The decomposition is then

$$\begin{aligned}
 & \left[ \begin{array}{c|c|c|c} \mathbf{Z}_{11} & \mathbf{Z}_{1,1} & \mathbf{X}_{11} & \mathbf{y}_{11} \\ \hline 1 & 1 & 0.8 & 1 & 0.8 & 0.8^2 & 0.92 \\ 1 & 1 & 1.2 & 1 & 1.2 & 1.2^2 & 3.91 \\ 1 & 1 & 1.6 & 1 & 1.6 & 1.6^2 & 7.69 \\ 1 & 1 & 2.0 & 1 & 2.0 & 2.0^2 & 11.80 \\ 1 & 1 & 2.4 & 1 & 2.4 & 2.4^2 & 16.00 \\ \hline \Delta_2 = 0.84 & 0 & 0 & 0 & 0 & 0 & 0 \end{array} \right] = \\
 & \mathbf{Q}_{(11)} \left[ \begin{array}{c|c|c|c} \mathbf{R}_{22(11)} & \mathbf{R}_{21(11)} & \mathbf{R}_{20(11)} & \mathbf{c}_{211} \\ \hline -2.389 & -2.093 & -3.349 & -2.093 & -3.349 & -6.029 & -16.880 \\ 0 & 0.087 & -0.024 & 0.087 & -0.024 & -0.528 & -1.343 \\ 0 & 0.087 & 0.376 & 0.087 & 0.376 & 0.592 & 2.437 \\ 0 & 0.087 & 0.776 & 0.087 & 0.776 & 2.032 & 6.547 \\ 0 & 0.087 & 1.176 & 0.087 & 1.176 & 3.792 & 10.747 \\ 0 & -0.767 & -1.029 & -0.767 & -1.029 & -1.653 & -4.412 \\ \hline 0 & \mathbf{R}_{11(11)} & \mathbf{R}_{10(11)} & \mathbf{c}_{1(11)} \end{array} \right]
 \end{aligned}$$

The direct contribution to the log-likelihood from this first site within the first wafer is

$$\log \text{abs} \left( \frac{|\Delta_2|}{|\mathbf{R}_{22(11)}|} \right) = \log (0.84/2.389) = -1.045$$

Since this was a balanced experiment where each wafer has the same number of sites and each site is measured at the same set of voltages, this contribution will be the same for each of the second level groups. The total direct contribution will be

$$\sum_{i=1}^{10} \sum_{j=1}^{67} \log (0.84/2.389) = -700.2$$

There will also be an indirect contribution from each site within each wafer according to the way they together determine  $\|\mathbf{c}_{-1}\|$  for a given value of  $\Delta_2$ .

## 2.5 Three or more levels of random effects

For every level of random effects added to the hierarchical model we simply extend the number of stages in the decompositions of the augmented  $\mathbf{Z}$  and  $\mathbf{X}$  matrices. For example, with three levels of random effects the decompositions begin with

$$\begin{aligned}
 & \left[ \begin{array}{c|c|c|c|c} \mathbf{Z}_{ijk} & \mathbf{Z}_{ij,k} & \mathbf{Z}_{i,jk} & \mathbf{X}_{ijk} & \mathbf{y}_{ijk} \\ \hline \Delta_3 & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} \end{array} \right] = \\
 & \mathbf{Q}_{(ijk)} \left[ \begin{array}{c|c|c|c|c} \mathbf{R}_{33(ijk)} & \mathbf{R}_{32(ijk)} & \mathbf{R}_{31(ijk)} & \mathbf{R}_{30(ijk)} & \mathbf{c}_{3(ijk)} \\ \hline \mathbf{0} & \mathbf{R}_{22(ijk)} & \mathbf{R}_{21(ijk)} & \mathbf{R}_{20(ijk)} & \mathbf{c}_{2(ijk)} \end{array} \right] \quad \begin{array}{l} i = 1, \dots, M \\ j = 1, \dots, M_i \\ k = 1, \dots, M_{ij} \end{array}
 \end{aligned}$$

to evaluate the integral associated with the third level of random effects. The next set of decompositions are

$$\begin{bmatrix} \mathbf{R}_{22(ij1)} & \mathbf{R}_{21(ij1)} & \mathbf{R}_{20(ij1)} & \mathbf{c}_{2(ij1)} \\ \vdots & \vdots & \vdots & \vdots \\ \mathbf{R}_{22(ijM_{ij})} & \mathbf{R}_{21(ijM_{ij})} & \mathbf{R}_{20(ijM_{ij})} & \mathbf{c}_{2(ijM_{ij})} \\ \mathbf{\Delta}_2 & \mathbf{0} & \mathbf{0} & \mathbf{0} \end{bmatrix} = \mathbf{Q}_{(ij)} \begin{bmatrix} \mathbf{R}_{22(ij)} & \mathbf{R}_{21(ij)} & \mathbf{R}_{20(ij)} & \mathbf{c}_{2(ij)} \\ \mathbf{0} & \mathbf{R}_{11(ij)} & \mathbf{R}_{10(ij)} & \mathbf{c}_{1(ij)} \end{bmatrix} \quad \begin{array}{l} i = 1, \dots, M \\ j = 1, \dots, M_i \\ k = 1, \dots, M_{ij} \end{array}$$

and from there we proceed as before with (13) and (7).

The profiled log-likelihood becomes

$$\begin{aligned} \ell(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \boldsymbol{\theta}_3 | \mathbf{y}) &= \log L(\widehat{\boldsymbol{\beta}}(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \boldsymbol{\theta}_3), \boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \boldsymbol{\theta}_3, \widehat{\sigma}^2(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \boldsymbol{\theta}_3) | \mathbf{y}) \\ &= \text{const} - N \log \|\mathbf{c}_{-1}\| + \sum_{i=1}^M \log \text{abs} \left( \frac{|\mathbf{\Delta}_1|}{|\mathbf{R}_{11(i)}|} \right) \\ &\quad + \sum_{i=1}^M \sum_{j=1}^{M_i} \log \text{abs} \left( \frac{|\mathbf{\Delta}_2|}{|\mathbf{R}_{22(ij)}|} \right) + \sum_{i=1}^M \sum_{j=1}^{M_i} \sum_{k=1}^{M_{ij}} \log \text{abs} \left( \frac{|\mathbf{\Delta}_3|}{|\mathbf{R}_{33(ijk)}|} \right) \end{aligned}$$

and the profiled restricted log-likelihood becomes

$$\begin{aligned} \ell_R(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \boldsymbol{\theta}_3 | \mathbf{y}) &= \log L_R(\widehat{\boldsymbol{\beta}}_R(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \boldsymbol{\theta}_3), \boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \boldsymbol{\theta}_3, \widehat{\sigma}^2_R(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \boldsymbol{\theta}_3) | \mathbf{y}) \\ &= \text{const} - (N - p) \log \|\mathbf{c}_{-1}\| - \log \text{abs} |\mathbf{R}_{00}| + \sum_{i=1}^M \log \text{abs} \left( \frac{|\mathbf{\Delta}_1|}{|\mathbf{R}_{11(i)}|} \right) \\ &\quad + \sum_{i=1}^M \sum_{j=1}^{M_i} \log \text{abs} \left( \frac{|\mathbf{\Delta}_2|}{|\mathbf{R}_{22(ij)}|} \right) + \sum_{i=1}^M \sum_{j=1}^{M_i} \sum_{k=1}^{M_{ij}} \log \text{abs} \left( \frac{|\mathbf{\Delta}_3|}{|\mathbf{R}_{33(ijk)}|} \right) \end{aligned}$$

## 2.6 Starting values for the $\boldsymbol{\theta}$ parameters

Because we can express both the profiled log-likelihood and the profiled log-restricted-likelihood as a function of the  $\boldsymbol{\theta}$  parameters, we only need to formulate starting values for  $\boldsymbol{\theta}$  when performing iterative optimization. From (10) and from the discussion of the example in §2.4, we can see that  $\boldsymbol{\theta}$  influences the profiled log-likelihood indirectly by changing  $\|\mathbf{c}_{-1}\|^2$  and directly through terms of the form

$$\begin{aligned} \log \text{abs} \left( \frac{|\mathbf{\Delta}|}{|\mathbf{R}_{11(i)}|} \right) &= -\frac{1}{2} \log \left( \frac{|\mathbf{Z}'_i \mathbf{Z}_i + \mathbf{\Delta}' \mathbf{\Delta}|}{|\mathbf{\Delta}' \mathbf{\Delta}|} \right) \\ &= \frac{-\log |\mathbf{I} + (\mathbf{\Delta}^{-1})' \mathbf{Z}'_i \mathbf{Z}_i \mathbf{\Delta}^{-1}|}{2} \leq 0 \end{aligned} \tag{14}$$

The size of  $\Delta$  relative to each of the  $Z_i, i = 1, \dots, M$  affects the two types of terms in opposite ways.

As  $\Delta$  gets larger it pulls the conditional least-squares values of  $b_i$  closer to 0 thus increasing the residual sum-of-squares term  $\|e_{-1}\|^2$ . Conversely when  $\Delta$  is very small relative to the  $Z_i$  it has little effect on the conditional least-squares values for  $b_i$  and  $\|e_{-1}\|^2$  will be smaller. However, small values of  $\Delta$  produce large values of the terms (14). For example, we showed earlier that  $\Delta_2 = 0.84$  produces

$$\log \text{abs} \left( \frac{|\Delta_2|}{|R_{22(11)}|} \right) = \log (0.84/2.389) = -1.045$$

in the IC manufacturing example. Doing the same calculations with a very small value, say  $\Delta_2 = 0.0084$  would produce

$$\log \text{abs} \left( \frac{|\Delta_2|}{|R_{22(11)}|} \right) = \log (0.0084/2.361) = -5.584$$

For a large value, say  $\Delta_2 = 84$ , we would get

$$\log \text{abs} \left( \frac{|\Delta_2|}{|R_{22(11)}|} \right) = \log (84/84.030) = -0.0003542$$

In the limit as  $\Delta$  becomes very large relative to the  $Z_i$ , the coefficients  $b_i$  in the penalized regression model are forced to zero and the residual sum-of-squares  $\|e_{-1}\|^2$  is the same as that from a regression of all the data on the fixed effects  $\beta$  only. The terms in (14) approach  $-\log |I|/2 = 0$ , which is their upper bound.

The net result is that extremely large values of  $\Delta$  or extremely small values of  $\Delta$  tend to produce small values of the likelihood for most data sets. The optimal values of  $\theta$  will usually correspond to a  $\Delta$  that is comparable in size to the  $Z_i$ . We generate a simple starting value  $\Delta^{(0)}$  as a diagonal matrix where each diagonal element is some fraction  $f$  of root-mean-square length of the corresponding column of the  $Z_i$  matrices. That is, letting  $Z_i(k)$  denote the  $k$ th column of  $Z_i$ , the initial value for the  $k$ th diagonal element of  $\Delta$  is  $f \times \left( \sum_{i=1}^M \|Z_i(k)\|^2 / M \right)^{1/2}$ . Some limited experimentation with this formulation indicated that fractions  $f$  between 1/4 and 1/2 worked well. We use  $f = 0.375$ .

In some patterned variance-covariance matrices, such as the compound symmetry matrices, a general diagonal matrix cannot be represented. For those cases we use as a starting estimate for  $\theta$  the parameter values that provide the patterned matrix that is closest to the diagonal  $\Delta$  calculated as described above.

This procedure generalizes easily to multiple levels of random effects. As an example, we consider the two-level model for the IC manufacturing data of §1.3. Because the design is balanced (i.e. all wafers have the same number of sites and the intensity of current measurements were made at the same voltages for every site), it suffices to calculate the column norms for, say,  $Z_{11}$  and  $Z_1$ . Therefore, the initial estimates for

the precision matrices are

$$\begin{aligned}\Delta_1^{(0)} &= 0.375 \begin{bmatrix} \|\mathbf{Z}_1(1)\| & 0 \\ 0 & \|\mathbf{Z}_1(2)\| \end{bmatrix} = 0.375 \begin{bmatrix} \sqrt{335} & 0 \\ 0 & \sqrt{964.8} \end{bmatrix} = \begin{bmatrix} 6.86 & 0 \\ 0 & 11.65 \end{bmatrix} \\ \Delta_2^{(0)} &= 0.375 \|\mathbf{Z}_{11}\| = 0.375\sqrt{5} = 0.84\end{aligned}$$

### 3 Conditional estimates and EM iterations

In the last section we derived a compact representation of the profiled log-likelihood or profiled log-restricted-likelihood for the single level or nested mixed-effects model. This representation used the relative precision factors  $\Delta_q$ ,  $q = 1, \dots, Q$  and orthogonal-triangular decompositions of augmented data arrays. A general nonlinear optimization algorithm, combined with the starting estimates for the  $\Delta_q$  derived in that section, could be used to obtain the maximum likelihood estimates,  $\hat{\boldsymbol{\theta}}$ , or the restricted maximum likelihood estimates,  $\hat{\boldsymbol{\theta}}_R$ , and, through (9),  $\hat{\boldsymbol{\beta}}$ .

Before considering details of the optimization we examine the general conditional estimates  $\hat{\boldsymbol{\beta}}(\boldsymbol{\theta})$  and the conditional distribution of the random effects  $p(\mathbf{b}_i | \mathbf{y}_i, \boldsymbol{\beta}, \boldsymbol{\theta}, \sigma^2)$  more closely. We relate these conditional estimates of the fixed-effects and the conditional expectations of the random effects to the parameter estimates in a penalized linear least squares problem. This conditional distribution of the random effects can be used to define EM iterations which are very helpful in refining starting estimates for the  $\Delta$ .

#### 3.1 Conditional distribution of the random effects

We can see from (5) that the conditional distribution of the random effects in a single-level model is

$$\begin{aligned}p(\mathbf{b}_i | \mathbf{y}_i, \boldsymbol{\beta}, \boldsymbol{\theta}, \sigma^2) &\propto p(\mathbf{b}_i, \mathbf{y}_i | \boldsymbol{\beta}, \boldsymbol{\theta}, \sigma^2) \\ &\propto \exp\left(\frac{-1}{2\sigma^2} \|\mathbf{c}_{1(i)} - \mathbf{R}_{10(i)}\boldsymbol{\beta} - \mathbf{R}_{11(i)}\mathbf{b}_i\|^2\right)\end{aligned}$$

or, writing all the random effects together as  $\mathbf{b}' = (\mathbf{b}'_1, \dots, \mathbf{b}'_M)'$ ,

$$p(\mathbf{b} | \mathbf{y}, \boldsymbol{\beta}, \boldsymbol{\theta}, \sigma^2) \propto \exp\left(\frac{-1}{2\sigma^2} \|\mathbf{c}_b - \mathbf{R}_b\mathbf{b} - \mathbf{R}_\beta\boldsymbol{\beta}\|^2\right)$$

where

$$\mathbf{R}_b = \begin{bmatrix} \mathbf{R}_{11(1)} & \mathbf{0} & \cdots & \mathbf{0} \\ \mathbf{0} & \mathbf{R}_{11(2)} & \cdots & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{0} & \mathbf{0} & \cdots & \mathbf{R}_{11(M)} \end{bmatrix} \quad \mathbf{R}_\beta = \begin{bmatrix} \mathbf{R}_{10(1)} \\ \mathbf{R}_{10(2)} \\ \vdots \\ \mathbf{R}_{10(M)} \end{bmatrix} \quad \mathbf{c}_b = \begin{bmatrix} \mathbf{c}_{1(1)} \\ \mathbf{c}_{1(2)} \\ \vdots \\ \mathbf{c}_{1(M)} \end{bmatrix}$$

It then follows that

$$\mathbf{b}|\mathbf{y} \sim \mathcal{N}\left(\mathbf{R}_b^{-1}(\mathbf{c}_b - \mathbf{R}_\beta\beta), \sigma^2 \mathbf{R}_b^{-1}(\mathbf{R}_b^{-1})'\right) \quad (15)$$

We note that  $\mathbf{R}_b^{-1}$  is a block diagonal matrix with its  $i$ th diagonal block given by the inverse of the  $q \times q$  upper-triangular matrix  $\mathbf{R}_{11(i)}$ .

From (8) we can see that the likelihood for the fixed effects  $\beta$  conditional on  $\theta$  and  $\sigma^2$  has a similar form

$$L(\beta|\mathbf{y}, \theta, \sigma^2) \propto \exp\left(\frac{-\|\mathbf{c}_0 - \mathbf{R}_{00}\beta\|^2}{2\sigma^2}\right)$$

The distribution of the conditional estimate  $\hat{\beta}(\theta)$  and the conditional distribution of  $\mathbf{b}$  are the same as the distribution of the estimate of  $\beta$  and the conditional estimates  $\mathbf{b}|\beta$  in a linear regression model

$$\mathbf{c} = \mathbf{R}\mathbf{b}^* + \boldsymbol{\epsilon}^* \quad \boldsymbol{\epsilon}^* \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$$

where

$$\mathbf{b}^* = \begin{bmatrix} \mathbf{b} \\ \beta \end{bmatrix} \quad \mathbf{R} = \begin{bmatrix} \mathbf{R}_b & \mathbf{R}_\beta \\ \mathbf{0} & \mathbf{R}_{00} \end{bmatrix} \quad \mathbf{c} = \begin{bmatrix} \mathbf{c}_b \\ \mathbf{c}_0 \end{bmatrix} \quad (16)$$

or, equivalently, the penalized linear regression model

$$\mathbf{y}^* = \mathbf{X}^* \mathbf{b}^* + \boldsymbol{\epsilon}^{**} \quad \boldsymbol{\epsilon}^{**} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$$

where

$$\mathbf{y}^* = \begin{bmatrix} \mathbf{y}_1 \\ \mathbf{0} \\ \mathbf{y}_2 \\ \mathbf{0} \\ \vdots \\ \mathbf{y}_M \\ \mathbf{0} \end{bmatrix} \quad \mathbf{X}^* = \begin{bmatrix} \mathbf{Z}_1 & \mathbf{0} & \cdots & \mathbf{0} & \mathbf{X}_1 \\ \boldsymbol{\Delta} & \mathbf{0} & \cdots & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{Z}_2 & \cdots & \mathbf{0} & \mathbf{X}_2 \\ \mathbf{0} & \boldsymbol{\Delta} & \cdots & \mathbf{0} & \mathbf{0} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{Z}_M & \mathbf{X}_M \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \boldsymbol{\Delta} & \mathbf{0} \end{bmatrix} \quad (17)$$

The conditional distribution of the random effects generalizes to multiple levels. For example, in a model with two levels of random effects,

$$\begin{aligned} p(\mathbf{b}_i, \mathbf{b}_{i1}, \dots, \mathbf{b}_{iM_i} | \mathbf{y}_{i1}, \dots, \mathbf{y}_{iM_i}, \beta, \theta, \sigma^2) &\propto \\ \exp \left[ \frac{-1}{2\sigma^2} \left( \|\mathbf{c}_{1(i)} - \mathbf{R}_{11(i)}\mathbf{b}_i - \mathbf{R}_{10(i)}\beta\|^2 + \right. \right. \\ &\quad \left. \left. \sum_{j=1}^{M_i} \|\mathbf{c}_{2(ij)} - \mathbf{R}_{22(ij)}\mathbf{b}_{ij} - \mathbf{R}_{21(ij)}\mathbf{b}_i - \mathbf{R}_{20(ij)}\beta\|^2 \right) \right] \end{aligned}$$

In a vector notation of  $\mathbf{b}(i)' = (\mathbf{b}'_{i1}, \dots, \mathbf{b}'_{iM_i}, \mathbf{b}'_i)'$  for the two-level model and writing

$$\mathbf{R}_b(i) = \begin{bmatrix} \mathbf{R}_{22(i1)} & \cdots & \mathbf{0} & \mathbf{R}_{21(i1)} \\ \vdots & \vdots & \vdots & \vdots \\ \mathbf{0} & \cdots & \mathbf{R}_{22(iM_i)} & \mathbf{R}_{21(iM_i)} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{R}_{11(i)} \end{bmatrix}$$

$$\mathbf{R}_\beta(i) = \begin{bmatrix} \mathbf{R}_{20(i1)} \\ \vdots \\ \mathbf{R}_{20(iM_i)} \\ \mathbf{R}_{10(i)} \end{bmatrix} \quad \mathbf{c}_b(i) = \begin{bmatrix} \mathbf{c}_{2(i1)} \\ \vdots \\ \mathbf{c}_{2(iM_i)} \\ \mathbf{c}_{1(i)} \end{bmatrix}$$

we obtain

$$\mathbf{b}(i)|\mathbf{y} \sim \mathcal{N} \left( \mathbf{R}_b^{-1}(i) (\mathbf{c}_b(i) - \mathbf{R}_\beta(i)\beta), \sigma^2 \mathbf{R}_b^{-1}(i) (\mathbf{R}_b^{-1}(i))' \right) \quad (18)$$

The equivalence of the distribution of the conditional estimates  $\hat{\beta}(\theta)$  and the conditional distribution of the random effects to the distribution of the parameter estimates in a penalized linear regression model extends to multiple levels of random effects. For two levels of random effects the penalized regression model is  $\mathbf{y}^* = \mathbf{X}^* \mathbf{b}^* + \epsilon^{**}$  where  $(\mathbf{b}^*)' = (\mathbf{b}(1)', \mathbf{b}(2)', \dots, \mathbf{b}(M)', \beta)'$  and the substitutions

$$\begin{bmatrix} \mathbf{Z}_i \\ \Delta \end{bmatrix} \Rightarrow \begin{bmatrix} \mathbf{Z}_{i1} & \mathbf{0} & \cdots & \mathbf{0} & \mathbf{Z}_{i,1} \\ \Delta_2 & \mathbf{0} & \cdots & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{Z}_{i2} & \cdots & \mathbf{0} & \mathbf{Z}_{i,2} \\ \mathbf{0} & \Delta_2 & \cdots & \mathbf{0} & \mathbf{0} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{Z}_{iM_i} & \mathbf{Z}_{i,M_i} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \Delta_2 & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \Delta_1 \end{bmatrix} \quad \begin{bmatrix} \mathbf{y}_i \\ \mathbf{0} \end{bmatrix} \Rightarrow \begin{bmatrix} \mathbf{y}_{i1} \\ \mathbf{0} \\ \mathbf{y}_{i2} \\ \mathbf{0} \\ \vdots \\ \mathbf{y}_{iM_i} \\ \mathbf{0} \\ \mathbf{0} \end{bmatrix}$$

are made in the form (17).

### 3.2 Conditional distributions and restricted likelihood

The expressions (15) and (18) and their generalizations to more levels of random effects allow the EM iterations on  $\theta$  for the likelihood criterion to be developed, as shown in §3.4. Similar expressions are used to develop EM iterations for the restricted likelihood criterion. As described in Laird and Ware (1982) the restricted likelihood can be regarded as the likelihood associated with model (1) but incorporating  $\beta$  as a random effect, say  $\mathbf{b}_0$ , with an associated relative variance matrix  $\mathbf{D}_0$  that tends to infinity. We could, for example, set  $\mathbf{D}_0 = k\mathbf{I}$ ,  $k \rightarrow \infty$ .

The conditional distribution of this extended random effects vector for the restricted model is simply the distribution of the estimates of  $\mathbf{b}^*$  in the corresponding regression models (16) or (17). That is,

$$\mathbf{b}^*|\mathbf{y} \sim \mathcal{N} (\mathbf{R}^{-1} \mathbf{c}, \sigma^2 \mathbf{R}^{-1} (\mathbf{R}^{-1})')$$



for a single level of random effects. This scales in the obvious way to multiple levels.

For both the original model and the restricted model the conditional distribution of the random effects corresponds to the distribution of the parameter estimates in a penalized linear regression model. The only difference is that in the original model it corresponds to the distribution of the estimates of  $\mathbf{b}$  conditional on  $\beta$  and in the restricted model it corresponds to the marginal distribution of the estimates of  $\mathbf{b}$ .

### 3.3 Using matrix sparsity

We have expressed the conditional distributions of the random effects in terms of the inverses of the upper triangular matrices  $\mathbf{R}$  and  $\mathbf{R}_b$ . These are large matrices but they are also patterned and sparse. We can take advantage of the patterns when forming the inverses and when doing further calculations.

For a single level of random effects we need only store

$$\begin{bmatrix} \mathbf{R}_{11(1)} & \mathbf{R}_{10(1)} & \mathbf{c}_{1(1)} \\ \mathbf{R}_{11(2)} & \mathbf{R}_{10(2)} & \mathbf{c}_{1(2)} \\ \vdots & \vdots & \vdots \\ \mathbf{R}_{11(M)} & \mathbf{R}_{10(M)} & \mathbf{c}_{1(M)} \\ \mathbf{0} & \mathbf{R}_{00} & \mathbf{c}_0 \end{bmatrix} \quad (19)$$

In calculating  $\mathbf{R}^{-1}$  we can take advantage of the fact that the pattern of zeroes in  $\mathbf{R}$  is repeated in

$$\mathbf{R}^{-1} = \begin{bmatrix} \mathbf{R}_{11(1)}^{-1} & \mathbf{0} & \cdots & \mathbf{0} & -\mathbf{R}_{11(1)}^{-1} \mathbf{R}_{10(1)} \mathbf{R}_{00}^{-1} \\ \mathbf{0} & \mathbf{R}_{11(2)}^{-1} & \cdots & \mathbf{0} & -\mathbf{R}_{11(2)}^{-1} \mathbf{R}_{10(2)} \mathbf{R}_{00}^{-1} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ \mathbf{0} & \mathbf{0} & \cdots & \mathbf{R}_{11(M)}^{-1} & -\mathbf{R}_{11(M)}^{-1} \mathbf{R}_{10(M)} \mathbf{R}_{00}^{-1} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{R}_{00}^{-1} \end{bmatrix} \quad (20)$$

and perform the operations in place on an array stored as in (19). That is, there is no need to expand the array  $\mathbf{R}$  out to its full size when calculating its inverse. We also note that each of the components such as  $\mathbf{R}_{00}$  and  $\mathbf{R}_{11(i)}$ ,  $i = 1, \dots, M$  that must be inverted to form (20) is a relatively small, triangular, non-singular matrix whose inverse is readily calculated.

One advantage of storing the intermediate results in this form is that the expressions for the conditional estimates  $\hat{\beta}(\theta)$  and  $\hat{\mathbf{b}}_i(\theta) = \mathbb{E}[\mathbf{b}_i | \mathbf{y}]$  evaluated at  $\hat{\beta}(\theta)$  have the form

$$\begin{aligned} \mathbf{R}_{00} \hat{\beta}(\theta) &= \mathbf{c}_0 \\ \mathbf{R}_{11(i)} \hat{\mathbf{b}}_i(\theta) &= \mathbf{c}_{1(i)} - \mathbf{R}_{10(i)} \hat{\beta}(\theta) & i = 1, \dots, M \\ \mathbf{R}_{22(ij)} \hat{\mathbf{b}}_{ij}(\theta) &= \mathbf{c}_{2(ij)} - \mathbf{R}_{20(ij)} \hat{\beta}(\theta) - \mathbf{R}_{21(ij)} \hat{\mathbf{b}}_i(\theta) & j = 1, \dots, M_i \\ &\dots & \dots \end{aligned}$$

Once  $\hat{\beta}(\theta)$  has been calculated, the evaluation of  $\mathbf{c}_{1(i)} - \mathbf{R}_{10(i)} \hat{\beta}(\theta)$ ,  $i = 1, \dots, M$  can be performed as a single matrix multiplication. The structure corresponding to (19)

for two levels of random effects has the form

$$\begin{bmatrix} \mathbf{R}_{22(11)} & \mathbf{R}_{21(11)} & \mathbf{R}_{20(11)} & \mathbf{c}_{2(11)} \\ \vdots & \vdots & \vdots & \vdots \\ \mathbf{R}_{22(1M_1)} & \mathbf{R}_{21(1M_1)} & \mathbf{R}_{20(1M_1)} & \mathbf{c}_{2(1M_1)} \\ \mathbf{0} & \mathbf{R}_{11(1)} & \mathbf{R}_{10(1)} & \mathbf{c}_{1(1)} \\ \mathbf{R}_{22(21)} & \mathbf{R}_{21(21)} & \mathbf{R}_{20(21)} & \mathbf{c}_{2(21)} \\ \vdots & \vdots & \vdots & \vdots \\ \mathbf{R}_{22(MM_M)} & \mathbf{R}_{21(MM_M)} & \mathbf{R}_{20(MM_M)} & \mathbf{c}_{2(MM_M)} \\ \mathbf{0} & \mathbf{R}_{11(M)} & \mathbf{R}_{10(M)} & \mathbf{c}_{1(M)} \\ \mathbf{0} & \mathbf{0} & \mathbf{R}_{00} & \mathbf{c}_0 \end{bmatrix} \quad (21)$$

Again, this form facilitates calculation of expressions such as  $\mathbf{c}_{1(i)} - \mathbf{R}_{10(i)}\hat{\boldsymbol{\beta}}(\boldsymbol{\theta})$  and  $\mathbf{c}_{2(ij)} - \mathbf{R}_{21(i)}\hat{\mathbf{b}}_i(\boldsymbol{\theta}) - \mathbf{R}_{20(ij)}\hat{\boldsymbol{\beta}}(\boldsymbol{\theta})$ .

The corresponding block in  $\mathbf{R}^{-1}$  is

$$\begin{bmatrix} \mathbf{R}_{22(i1)}^{-1} & \mathbf{0} & \cdots & \mathbf{0} & -\mathbf{R}_{22(i1)}^{-1}\mathbf{R}_{21(i1)}\mathbf{R}_{11(1)}^{-1} \\ \mathbf{0} & \mathbf{R}_{22(i2)}^{-1} & \cdots & \mathbf{0} & -\mathbf{R}_{22(i2)}^{-1}\mathbf{R}_{21(i2)}\mathbf{R}_{11(1)}^{-1} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ \mathbf{0} & \mathbf{0} & \cdots & \mathbf{R}_{22(iM_i)}^{-1} & -\mathbf{R}_{22(iM_i)}^{-1}\mathbf{R}_{21(iM_i)}\mathbf{R}_{11(i)}^{-1} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{R}_{11(i)}^{-1} \end{bmatrix} \quad (22)$$

We can see that the pattern of (20) is repeated here. As before, the block does not need to be expanded to its full size to calculate the components of the inverse—the condensed array of the form (21) can be manipulated in place.

### 3.4 EM iterations

Optimization of the profiled log-likelihood or restricted log-likelihood is usually accomplished through EM iterations or through Newton-Raphson iterations (Laird and Ware, 1982; Lindstrom and Bates, 1988; Longford, 1993). The EM iterations have the advantage that the individual iterations are easy to compute and the initial iterations approach the optimum quite quickly. However, close to the optimum the EM iterations often proceed very slowly. It can be difficult to decide if the EM iterations have converged when they end up taking an exceedingly large number of very small steps toward the optimum. Also, taking so many iterations is itself expensive even if each iteration is relatively simple and fast.

The Newton-Raphson iterations, on the other hand, are individually more computationally intensive than the EM iterations and they can be quite unstable when far from the optimum. However, close to the optimum they converge very quickly.

We therefore recommend a hybrid approach of forming an initial  $\boldsymbol{\theta}^{(0)}$  as described in section 2.6, performing a moderate number of EM iterations, then switching to Newton-Raphson iterations. Essentially the EM iterations can be regarded as refining the starting estimates before beginning the more general optimization routine.

The EM iterations are based on regarding the random effects, such as the  $\mathbf{b}_i$ ,  $i = 1, \dots, M$ , as unobserved data. At iteration  $t$  we use the current variance-covariance parameter vector,  $\boldsymbol{\theta}^{(t)}$ , to evaluate the distribution of  $\mathbf{b}|\mathbf{y}$  and derive the expectation of the log-likelihood for a new value of  $\boldsymbol{\theta}$  given this conditional distribution. Because we are taking an expectation, this step is called the E step.

The M step consists of maximizing this expectation with respect to  $\boldsymbol{\theta}$  to produce  $\boldsymbol{\theta}^{(t+1)}$ .

We can perform the expectation and maximization steps with respect to  $\boldsymbol{\theta}$  only because we have explicit forms for the profiled likelihood as a function of  $\boldsymbol{\theta}$ . At iteration  $t$  we derive  $\boldsymbol{\theta}^{(t+1)}$  from  $\boldsymbol{\theta}^{(t)}$  fixing  $\boldsymbol{\beta} = \widehat{\boldsymbol{\beta}}(\boldsymbol{\theta}^{(t)})$  and  $\sigma^2 = \widehat{\sigma}^2(\boldsymbol{\theta}^{(t)})$ . The EM algorithm ensures that the log-likelihood will increase from  $\boldsymbol{\theta}^{(t)}$  to  $\boldsymbol{\theta}^{(t+1)}$ . If we then adjust  $\boldsymbol{\beta}$  and  $\sigma^2$  to the optimal values given  $\boldsymbol{\theta}^{(t+1)}$  we can only increase the log-likelihood further.

For a single level of random effects we can see from (3) that the likelihood for the “full data”—both the unobserved  $\mathbf{b}$  and the observed  $\mathbf{y}$ —depends on  $\boldsymbol{\theta}$  only through  $|\mathbf{D}|$  and quadratic forms in  $(\mathbf{b}_i/\sigma)' \mathbf{D}^{-1} \mathbf{b}_i/\sigma$ . To obtain the expected value for a new value of  $\boldsymbol{\theta}$  we use standard results on the distribution of quadratic forms plus the expressions in (15) to derive

$$\begin{aligned} \mathbb{E}[(\mathbf{b}_i/\sigma)' \mathbf{D}_1^{-1} \mathbf{b}_i/\sigma | \mathbf{y}] &= \mathbb{E}[\mathbf{b}_i/\sigma | \mathbf{y}]' \mathbf{D}_1^{-1} \mathbb{E}[\mathbf{b}_i/\sigma | \mathbf{y}] + \text{tr}(\mathbf{D}^{-1} \text{Var}(\mathbf{b}_i | \mathbf{y})) \\ &= \|\boldsymbol{\Delta}_1 \mathbb{E}[\mathbf{b}_i | \mathbf{y}]/\sigma\|^2 + \left\| \left( \mathbf{R}_{11(i)}^{-1} \right)' \boldsymbol{\Delta}_1' \right\|_2^2 \\ &= \left\| \begin{bmatrix} \mathbb{E}[\mathbf{b}_i | \mathbf{y}]'/\sigma \\ \left( \mathbf{R}_{11(i)}^{-1} \right)' \boldsymbol{\Delta}_1' \end{bmatrix} \right\|_2^2 \end{aligned}$$

where the 2-norm of a matrix, written  $\|\mathbf{A}\|_2$ , is the square root of the sum of squares of the entries in the matrix. In this expression  $\mathbf{R}_{11(i)}$  is evaluated using the data and  $\boldsymbol{\theta}^{(t)}$  because it is a characteristic of the conditional distribution but  $\boldsymbol{\Delta}_1$  is evaluated at the general  $\boldsymbol{\theta}$  over which we will be optimizing the expected likelihood given the conditional distribution so as to produce  $\boldsymbol{\theta}^{(t+1)}$ .

The sum of the contributions for the level-1 random effects can be expressed as

$$\left\| \begin{bmatrix} \mathbb{E}[\mathbf{b}_1 | \mathbf{y}]'/\sigma \\ \left( \mathbf{R}_{11(1)}^{-1} \right)' \\ \vdots \\ \mathbb{E}[\mathbf{b}_M | \mathbf{y}]'/\sigma \\ \left( \mathbf{R}_{11(M)}^{-1} \right)' \end{bmatrix} \boldsymbol{\Delta}_1' \right\|_2^2 = \|\mathbf{A}_1 \boldsymbol{\Delta}_1'\|_2^2 \quad (23)$$

where  $\mathbf{U}_1 \mathbf{A}_1$  is an orthogonal-triangular decomposition of the stacked matrix on the left side of (23). Thus  $\mathbf{U}_1$  is a  $M(q_1 + 1) \times q_1$  matrix with orthonormal columns and  $\mathbf{A}_1$  is a  $q_1 \times q_1$  upper-triangular, non-singular matrix.

As described in §3.1, the conditional distribution of  $\mathbf{b}|\mathbf{y}$  for the restricted model is also multivariate normal with the same mean but with a different variance-covariance matrix. For this case the non-zero blocks of the form  $-\mathbf{R}_{11(1)}^{-1} \mathbf{R}_{10(1)} \mathbf{R}_{00}^{-1}$  in the  $\mathbf{R}^{-1}$

matrix shown in (20) must be incorporated in the calculation of  $\mathbf{A}_{R1}$ . The decomposition has the form

$$\begin{bmatrix} \mathbf{E}[\mathbf{b}_1|\mathbf{y}]'/\sigma \\ \left(\mathbf{R}_{11(1)}^{-1}\right)' \\ -\left(\mathbf{R}_{11(1)}^{-1}\mathbf{R}_{10(1)}\mathbf{R}_{00}^{-1}\right)' \\ \vdots \\ \mathbf{E}[\mathbf{b}_M|\mathbf{y}]'/\sigma \\ \left(\mathbf{R}_{11(M)}^{-1}\right)' \\ -\left(\mathbf{R}_{11(M)}\mathbf{R}_{10(M)}\mathbf{R}_{00}^{-1}\right)' \end{bmatrix} = \mathbf{U}_{R1}\mathbf{A}_{R1}$$

With two levels of random effects, the matrix  $\mathbf{A}_1$  is calculated as in (23) and the matrix  $\mathbf{A}_2$  is calculated from a decomposition of a matrix formed from the scaled estimates of the random effects at that level and their estimated variances calculated according to (22).

$$\begin{bmatrix} \mathbf{E}[\mathbf{b}_{11}|\mathbf{y}]'/\sigma \\ \left(\mathbf{R}_{22(11)}^{-1}\right)' \\ -\left(\mathbf{R}_{22(11)}^{-1}\mathbf{R}_{21(11)}\mathbf{R}_{11(1)}^{-1}\right)' \\ \vdots \\ \mathbf{E}[\mathbf{b}_{MM_M}|\mathbf{y}]'/\sigma \\ \left(\mathbf{R}_{22(MM_M)}^{-1}\right)' \\ -\left(\mathbf{R}_{22(MM_M)}^{-1}\mathbf{R}_{21(MM_M)}\mathbf{R}_{11(M)}^{-1}\right)' \end{bmatrix} = \mathbf{U}_2\mathbf{A}_2$$

In general,  $\mathbf{A}_j$  is calculated by computing the expected values of the conditional distribution of the level  $j$  random effects and scaling them by  $1/\sigma$ . The non-zero blocks in  $\mathbf{R}_b^{-1}$  are calculated and transposed. All those blocks corresponding to the level- $j$  random effects are stacked, along with the scaled expected random effects. The resulting matrix of  $q_j$  columns is decomposed to give the upper-triangular matrix  $\mathbf{A}_j$ . When computing  $\mathbf{A}_{Rj}$  for iterations based on the restricted likelihood, the matrix  $\mathbf{R}_b^{-1}$  is replaced by  $\mathbf{R}^{-1}$ .

Once  $\mathbf{A}_1$  (or  $\mathbf{A}_{R1}$ ) is available, the maximization problem to determine  $\theta_1^{(t+1)}$  reduces to

$$\max_{\theta_1} \{M \log |\mathbf{D}_1^{-1}| - \text{tr}(\mathbf{A}_1' \mathbf{D}_1^{-1} \mathbf{A}_1)\} \quad (24)$$

If  $\mathbf{D}_1$  can be a general positive-definite matrix, the maximizer of (24) produces  $\mathbf{D}_1^{-1} = \Delta_1' \Delta_1$  where

$$\Delta_1^{(t+1)} = \frac{(\mathbf{A}_1^{-1})'}{\sqrt{M}} \quad (25)$$

Since only the factor  $\Delta_1$  and not the corresponding  $\theta$  is needed in the calculation of the next iteration, the iteration loop is complete. For two levels of random effects, the optimization problem becomes

$$\max_{\theta_2} \left\{ \left( \sum_{i=1}^M M_i \right) \log |D_2^{-1}| - \text{tr}(A_2' D_2^{-1} A_2) \right\}$$

with a solution, for general positive-definite  $D_2$  of

$$\Delta_2^{(t+1)} = \frac{(A_2^{-1})'}{\sqrt{\sum_{i=1}^M M_i}} \quad (26)$$

This scales in the obvious way.

The only change required to perform EM iterations for the restricted likelihood criterion is to replace  $A_j$  by  $A_{Rj}$  in (25) or (26) or their extended versions for more levels.

In the IC manufacturing example of §1.3, starting with the initial estimates  $\Delta_1^{(0)}$  and  $\Delta_2^{(0)}$  given in §2.6, we obtain after 20 EM iterations

$$\Delta_1^{(20)} = \begin{bmatrix} 2.498 & 0 \\ 1.097 & 0.681 \end{bmatrix} \quad \Delta_2^{(20)} = 0.477$$

The maximum likelihood estimates for this example are

$$\hat{\Delta}_1 = \begin{bmatrix} 2.501 & 0 \\ 1.091 & 0.681 \end{bmatrix} \quad \hat{\Delta}_2 = 0.477$$

indicating that the EM algorithm quickly placed the estimates of  $\theta$  in the neighborhood of the MLEs.

### 3.5 General optimization and obtaining BLUPs

After refining the starting estimates for  $\theta$  with a moderate number of EM iterations, the parameters  $\theta$  are usually close to their optimal values. A general unconstrained optimization algorithm such as Newton-Raphson or a quasi-Newton algorithm can then be used to produce the final parameter estimates  $\hat{\theta}$ . The maximum likelihood estimates for  $\beta$  can be evaluated from  $R_{00}$  and  $c_0$  at  $\hat{\theta}$  and (7). Furthermore, the matrices  $R_{11(i)}$  and  $R_{10(i)}$  and the vectors  $c_{1(i)}$  can be used with  $\hat{\beta}$  to solve for  $E[b_i|y]$  at the parameter estimates. These are called the Best Linear Unbiased Predictors or BLUPs for the random effects.

When performing the optimization it is often helpful if the gradient of the objective function and possibly the Hessian of the objective function can be evaluated. It may be possible to use techniques similar to those in Golub and Pereyra (1973) or Kaufman (1975) and in Bates (1983) to derive an expression for the gradient using the penalized linear regression representation of the model. This will be difficult to generalize to several levels of nested random effects. Because the profiled log-likelihood or profiled log-restricted-likelihood can be evaluated efficiently for different values of  $\theta$ , we prefer to use numerical values of the gradient and Hessian calculated through finite differences.

## 4 Computational considerations

### 4.1 Organizing the decompositions

Although several different matrices are being calculated and decomposed in the process of evaluating the profiled log-likelihood for  $\theta$ , most of these matrices do not need to be stored unless  $\hat{\beta}$  or the BLUPs for the  $b_i$  are to be calculated. In fact, the total amount of storage required for the calculation is essentially the same as that required to represent the original data.

The decompositions can be performed “in place” with a slightly augmented version of the original data array or they can be performed by copying horizontal slices of the original data array to temporary storage, augmenting them, decomposing the result, then copying back the pieces that are to be used later.

To demonstrate these techniques we begin with the data for a model with a single level of random effects, organized as

$$\begin{bmatrix} \mathbf{Z}_1 & \mathbf{X}_1 & \mathbf{y}_1 \\ \vdots & \vdots & \vdots \\ \mathbf{Z}_M & \mathbf{X}_M & \mathbf{y}_M \end{bmatrix}$$

For a given value of  $\theta$  we evaluate  $\Delta$  then, for each of the  $M$  horizontal slices in the array, form and decompose

$$\begin{bmatrix} \mathbf{Z}_i & \mathbf{X}_i & \mathbf{y}_i \\ \Delta & \mathbf{0} & \mathbf{0} \end{bmatrix} = \mathbf{Q}_{(i)} \begin{bmatrix} \mathbf{R}_{11(i)} & \mathbf{R}_{10(i)} & \mathbf{c}_{1(i)} \\ \mathbf{0} & \mathbf{R}_{00(i)} & \mathbf{c}_{0(i)} \end{bmatrix}$$

exactly as in (4). The determinant of  $\mathbf{R}_{11(i)}$  is calculated as the product of its diagonal elements and stored. After that none of  $\mathbf{R}_{11(i)}$ ,  $\mathbf{R}_{10(i)}$ , or  $\mathbf{c}_{1(i)}$  are required for the evaluation of the log-likelihood. They are only needed if the BLUP for  $b_i$  is to be calculated as the solution to  $\mathbf{R}_{11(i)}\mathbf{E}[b_i|\mathbf{y}] = \mathbf{c}_{1(i)} - \mathbf{R}_{10(i)}\hat{\beta}$ . After the decomposition, the matrix  $\mathbf{R}_{00(i)}$  and the vector  $\mathbf{c}_{0(i)}$  can be copied back into the storage previously occupied by  $\mathbf{X}_i$  and  $\mathbf{y}_i$ . When we have finished doing this for  $i = 1, \dots, M$  the last  $p + 1$  columns of the original data array will be in exactly the form needed to evaluate the decomposition (7). This provides  $\mathbf{c}_{-1}$  and  $\mathbf{R}_{00}$ . All the information to evaluate the profiled log-likelihood or the profiled log-restricted-likelihood is then available.

The operations of copying slices of the original data array to temporary storage then copying some of the results back into the data array can be avoided, if desired. Beginning with an augmented form of the data array as

$$\begin{bmatrix} \mathbf{Z}_1 & \mathbf{X}_1 & \mathbf{y}_1 \\ \vdots & \vdots & \vdots \\ \mathbf{Z}_M & \mathbf{X}_M & \mathbf{y}_M \\ \Delta & \mathbf{0} & \mathbf{0} \end{bmatrix}$$

we can decompose the augmented array for the  $M$ th case in place and evaluate  $\text{abs}|\mathbf{R}_{11(M)}|$ . The rows containing  $\mathbf{R}_{11(M)}$ ,  $\mathbf{R}_{10(M)}$ , and  $\mathbf{c}_{1(M)}$  can then be overwritten with  $\Delta$  in the

$Z$  columns and zeroes in the  $X$  and  $y$  columns. This results in the augmented arrays for the  $(M - 1)$ 'st case being available to be decomposed. We continue in this fashion working from the  $M$ 'th case up to the first case. The arrays needed to calculate the decomposition (7) are then in place in the last  $p + 1$  columns but beginning at row  $q + 1$ , not at the first row.

In our experience the cost of copying the arrays to temporary storage before decomposing and copying some results back after decomposing is negligible compared to the cost of decomposing the arrays. Thus we prefer to work on copies. Furthermore, even though the “in place” calculation does not require copying the data it may at times take longer than the “copying” calculation. This is because the memory references are widely scattered in the memory space when decomposing a horizontal slice of a matrix with a very large number of rows. When working with extremely large data sets that cannot fit into the available physical memory of the computer such access patterns will degrade the performance of virtual memory systems. In the approach based on copying the widely scattered memory locations are only used twice; once as the source for the copying and once as the destination of the copying. When actually doing the decomposition on the copy, the memory references are more localized.

With two levels of random effects we begin with the data array in the form

$$\begin{bmatrix} Z_{11} & Z_{1,1} & X_{11} & y_{11} \\ \vdots & \vdots & \vdots & \vdots \\ Z_{1M_1} & Z_{1,M_1} & X_{1M_1} & y_{1M_1} \\ Z_{21} & Z_{2,1} & X_{21} & y_{21} \\ \vdots & \vdots & \vdots & \vdots \\ Z_{MM_M} & Z_{M,M_M} & X_{MM_M} & y_{MM_M} \end{bmatrix}$$

and decompose the augmented matrices from each horizontal slice corresponding to the level-2 random effects while accumulating  $\text{abs} |R_{22(ij)}|$ ,  $i = 1, \dots, M$ ,  $j = 1, \dots, M_i$ . The matrices  $Z_{i,j}$ ,  $X_{ij}$ , and the vector  $y_{ij}$  are overwritten by  $R_{11(ij)}$ ,  $R_{10(ij)}$ , and  $c_{1(ij)}$  produced in this decomposition. The modified data array is then in the form needed to perform the decompositions for the level-1 random effects (ignoring the first  $q_2$  columns).

## 4.2 Pre-decomposition of the original data array

When the number of observations per lowest-level group is large compared to the number of random effects for the group, some time can be saved by decomposing the original data arrays and saving only the triangular parts of the decomposition. This technique is similar to the “two-stage orthogonal factorization” described by Golub and Pereyra (1973).

With one level of random effects the initial decomposition is similar to the reductions in (4) and (7). First we decompose

$$\begin{bmatrix} Z_i & X_i & y_i \end{bmatrix} = S_{(i)} \begin{bmatrix} T_{11(i)} & T_{10(i)} & d_{1(i)} \\ \mathbf{0} & T_{00(i)} & d_{0(i)} \end{bmatrix} \quad i = 1, \dots, M$$

then we accumulate and decompose all the  $T_{00(i)}$  matrices and  $c_{0(i)}$  vectors as

$$\begin{bmatrix} T_{00(1)} & d_{0(1)} \\ \vdots & \vdots \\ T_{00(M)} & d_{0(M)} \end{bmatrix} = S_0 \begin{bmatrix} T_{00} & d_0 \\ \mathbf{0} & d_{-1} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \quad (27)$$

where  $d_{-1}$  is a scalar. The decomposition (27) is the one occasion where we carry out the orthogonal-triangular decomposition completely to reduce the matrix to triangular form.

The information in

$$\begin{bmatrix} T_{11(1)} & T_{10(1)} & d_{1(1)} \\ \vdots & \vdots & \vdots \\ T_{11(M)} & T_{10(M)} & d_{1(M)} \\ \mathbf{0} & T_{00} & d_0 \\ \mathbf{0} & \mathbf{0} & d_{-1} \end{bmatrix}$$

can now be used in place of the original data arrays. The number of rows in each group is reduced to  $q$  except for the  $M$ 'th group where we include the rows with  $T_{00}$ ,  $d_0$ , and  $d_{-1}$ . These rows must be included in the overall data array but could be included with any of the groups. We include them with the last group for convenience and to emulate the pattern in (19).

With two levels of random effects the initial decompositions are of the form

$$\begin{bmatrix} Z_{ij} & Z_{i,j} & X_{ij} & y_{ij} \end{bmatrix} = S_{(ij)} \begin{bmatrix} T_{22(ij)} & T_{21(ij)} & T_{20(ij)} & d_{2(ij)} \\ \mathbf{0} & T_{11(ij)} & T_{10(ij)} & d_{1(ij)} \end{bmatrix} \quad \begin{array}{l} i = 1, \dots, M \\ j = 1, \dots, M_i \end{array}$$

followed by

$$\begin{bmatrix} T_{11(i1)} & T_{10(i1)} & d_{1(i1)} \\ \vdots & \vdots & \vdots \\ T_{11(iM_i)} & T_{10(iM_i)} & d_{1(iM_i)} \end{bmatrix} = S_{(i)} \begin{bmatrix} T_{11(i)} & T_{10(i)} & d_{1(i)} \\ \mathbf{0} & T_{00(i)} & d_{0(i)} \end{bmatrix} \quad i = 1, \dots, M$$

and finally (27). The information required for later calculations is stored in the same form as that in (21) with a  $R$  matrix replaced by the corresponding  $T$  matrix and a  $c$  vector replaced by the corresponding  $d$  vector.

Generalization to more than two levels of random effects follows this type of pattern. The similarity between the decompositions needed for the evaluation of the log-likelihood or restricted log-likelihood and those needed for the pre-decomposition step can be exploited when writing software. In our code we use a single function or method for the decomposition step with an additional argument for the matrix  $\Delta$  that should be appended as a new set of rows before decomposing. When doing the pre-decomposition the number of rows to be appended is set to zero. Because the information to be stored for calculation of  $\hat{\beta}$  and the BLUPs for the random effects or for the pre-decomposition is of a similar shape, we use the same function for both storage operations.



### 4.3 Parallelization

An important consideration in computational methods for modern computer systems is determining which parts of the calculation can be performed in parallel. The structure of the methods we have described divides neatly into separately evaluated pieces.

Within each level of random effects, the reductions and other computations related to the different groups are distinct. Thus they can be performed in parallel. Because all results from a given level of random effects must be available before proceeding to the next level, the extent of possible parallelization is clearly defined.

## 5 Conclusions and extensions to other models

Computational methods for maximum likelihood or restricted maximum likelihood estimation of the parameters in a linear multilevel mixed-effects model are greatly enhanced by expressing the variance-covariance matrix of the random effects at each level in terms of a square root of the inverse of the relative variance matrix. These are the matrices that we have written as  $\Delta_q, q = 1, \dots, Q$ .

Using this formulation and taking matrix decompositions, the profiled log-likelihood or profiled restricted log-likelihood can be compactly expressed and calculated. These expressions also give an indication of suitable starting values for the variance-covariance parameters.

The calculation of conditional estimates of the fixed effects or BLUPs for the random effects is also expressed compactly using matrix decomposition techniques. Using these expressions an EM algorithm for parameter estimation can be readily derived. The combination of starting estimates calculated from the original design matrices and a moderate number of EM iterations usually puts the parameter values very close to the final parameter estimates. The optimization can be finished with a few Newton-Raphson iterations in a suitable parameterization. Again, we can take advantage of the profiling of the log-likelihood over the values of the fixed-effects parameters.

The computational methods described in the previous sections can also be applied to extensions of the basic Gaussian linear multilevel model of §1. Two such extensions will be considered in this section: Gaussian linear multilevel models with *non-spherical* distributions for the within-group errors and Gaussian *nonlinear* multilevel models.

### 5.1 General linear multilevel model

The basic Gaussian linear multilevel model assumes that the within-group errors  $\epsilon_i$  are distributed as  $\mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$ . In many applications, especially when longitudinal or spatial data are collected, it is reasonable to allow for correlation among the within-group errors. The assumption of equal variances for the within-group errors also is frequently violated in practice. A more general formulation of the Gaussian multilevel model allows non-spherical Gaussian distributions for the within-group error.

In the single-level linear mixed-effects model (1), the general formulation of the model allows  $\epsilon_i \sim \mathcal{N}(\mathbf{0}, \sigma^2 \Lambda_i)$ , where the scaled variance-covariance matrix  $\Lambda_i$  depends on  $i$  only through its dimensions and is generally parametrized by a small, fixed

set of parameters  $\rho$ . We assume that the  $\Lambda_i$ ,  $i = 1, \dots, M$  are positive definite. This formulation allows both heteroscedastic (e.g. variance increasing with a power of the expected response) and correlated (e.g. autoregressive-moving average structures) within-group errors. For the general linear two-level model (2) we have  $\epsilon_{ij} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \Lambda_{ij})$  and this is similarly extended for more levels of nesting.

Let  $\Lambda_i^{1/2}$  denote a square-root factor of  $\Lambda_i$  (that is,  $\Lambda_i = (\Lambda_i^{1/2})' \Lambda_i^{1/2}$ ) and define

$$\mathbf{y}_i^* = (\Lambda_i^{-1/2})' \mathbf{y}_i \quad \mathbf{X}_i^* = (\Lambda_i^{-1/2})' \mathbf{X}_i \quad \mathbf{Z}_i^* = (\Lambda_i^{-1/2})' \mathbf{Z}_i \quad \epsilon_i^* = (\Lambda_i^{-1/2})' \epsilon_i$$

It then follows from elementary properties of the multivariate normal distribution that

$$\begin{aligned} \mathbf{y}_i^* &= \mathbf{X}_i^* \boldsymbol{\beta} + \mathbf{Z}_i^* \mathbf{b}_i + \epsilon_i^*, \quad i = 1, \dots, M \\ \mathbf{b}_i &\sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}), \quad \epsilon_i^* \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}) \end{aligned} \quad (28)$$

That is,  $\mathbf{y}_i^*$  follows the basic single-level linear mixed-effects (1). The Jacobian of the linear transformation  $\mathbf{y}_i \rightarrow \mathbf{y}_i^*$  is  $1/\text{abs}|\Lambda_i^{1/2}|$  and by (10) and (28)

$$\begin{aligned} \ell(\boldsymbol{\theta}, \boldsymbol{\rho} | \mathbf{y}) &= \text{const} - N \log \|\mathbf{c}_{-1}^*\| + \sum_{i=1}^M \log \text{abs} \left( \frac{|\Delta|}{|\mathbf{R}_{11(i)}^*|} \right) - \sum_{i=1}^M \log \text{abs} |\Lambda_i^{1/2}| \\ &= \ell(\boldsymbol{\theta} | \mathbf{y}^*) - \sum_{i=1}^M \log \text{abs} |\Lambda_i^{1/2}| \end{aligned}$$

where  $\mathbf{c}_{-1}^*$  and  $\mathbf{R}_{11(i)}^*$  are the equivalent of  $\mathbf{c}_{-1}$  and  $\mathbf{R}_{11(i)}$  for model (28). The restricted log-likelihood is given by

$$\ell_R(\boldsymbol{\theta}, \boldsymbol{\rho} | \mathbf{y}) = \ell_R(\boldsymbol{\theta} | \mathbf{y}^*) - \sum_{i=1}^M \log \text{abs} |\Lambda_i^{1/2}|$$

Similarly, the log-likelihood for the general linear two-level model is

$$\ell(\boldsymbol{\theta}, \boldsymbol{\rho} | \mathbf{y}) = \ell(\boldsymbol{\theta} | \mathbf{y}^*) - \sum_{i=1}^M \sum_{j=1}^{M_i} \log \text{abs} |\Lambda_{ij}^{1/2}|$$

with an equivalent expression for the restricted log-likelihood. This extends to an arbitrary number of levels in the obvious way.

The computational methods described in §4 can, for the most part, be easily extended to the general linear multilevel model by replacing  $\mathbf{y}$ ,  $\mathbf{X}$ , and  $\mathbf{Z}$  with  $\mathbf{y}^*$ ,  $\mathbf{X}^*$ , and  $\mathbf{Z}^*$ . For example, given an initial value  $\boldsymbol{\rho}^{(0)}$  for the within-group covariance structure parameters, initial values for  $\boldsymbol{\theta}$  can be obtained using the methodology of §2.6 (e.g.

$$\left[ \Delta^{(0)} \right]_{kk} = 0.375 \sqrt{\sum_{i=1}^M \|\mathbf{Z}_i^*(k)\|^2 / M} \text{ in the single-level case}.$$

The derivation of  $\hat{\beta}$  and the BLUPs for the general linear multilevel model is done as in §3.5, with the obvious substitutions. Pre-decompositions of the data, as described in §4.2, are not meaningful for the general linear multilevel model, as a new decomposition would have to be obtained for each new value of  $\rho$  in the optimization process. The results on the distributions of  $\hat{\beta}(\theta)$  and  $b|y$  in §3.1 extend in the obvious way to the general linear multilevel model.

Finally, the EM algorithm described in §3.4 has to be modified for the general linear multilevel model. For a given value of  $\rho = \rho^{(t)}$ , the methodology of §3.4 can be applied to the corresponding  $y^*$  “observations” to obtain updated estimates  $\theta^{(t)}$ ,  $\beta^{(t)}$ , and  $\sigma^{2(t)}$ . Assuming these fixed, updated estimates of  $\rho$  are obtained by maximizing  $\ell(\rho|y, \theta^{(t)}, \beta^{(t)}, \sigma^{2(t)})$ . This alternating optimization scheme is an example of the ECME algorithm proposed by Liu and Rubin (1994) and shares with the EM algorithm the property of monotone convergence.

## 5.2 Nonlinear multilevel model

A one-level nonlinear mixed effects model is similar in form to the linear mixed effects models (1) except that the expression  $X_i\beta + Z_ib_i$ , which is linear in both the fixed effects  $\beta$  and the random effects  $b_i$ , is replaced by a nonlinear expression  $f_i(\beta, b_i)$  where the components of  $f_i(\beta, b_i)$  are given by  $\{f_i(\beta, b_i)\}_j = f(\phi_{ij}, x_{ij})$ . Here  $f$  is a nonlinear model function and the subject-specific model parameter for group  $i$  at the  $j$ th observation is

$$\phi_{ij} = A_{ij}\beta + B_{ij}b_i \quad .$$

The matrices  $A_{ij}$  and  $B_{ij}$  are of appropriate dimension and depend on the group and possibly on the values of some covariates at the  $j$ th observation. This model is a slight generalization of that described in Lindstrom and Bates (1990) or Davidian and Giltinan (1995) in that  $A_{ij}$  and  $B_{ij}$  can depend on  $j$ . This generalization allows the incorporation of “time-varying” covariates with the fixed effects or the random effects in the model.

A nonlinear model with two levels of random effects can be written

$$y_{ij} = f_{ij}(\beta, b_i, b_{ij}) + \epsilon_{ij} \quad i = 1, \dots, M \quad j = 1, \dots, M_i$$

where  $\{f_{ij}(\beta, b_i, b_{ij})\}_k = f(\phi_{ijk}, x_{ijk})$  and  $\phi_{ijk} = A_{ijk}\beta + B_{i,j,k}b_i + B_{ij,k}b_{ij}$ . The reason for the somewhat awkward notation is, as before, to allow “time-varying” covariates to be used with the fixed effects or the random effects. Often in practice the matrices  $A_{ijk}$ ,  $B_{i,j,k}$ , and  $B_{ij,k}$  are identity matrices or subsets of the columns of an identity matrix.

Extensions to more than two levels of random effects follow the obvious patterns.

Because  $f$  can be nonlinear in the fixed or random effects, the integrals required to express the log-likelihood do not have the succinct representations of those in the linear case. Also the conditional estimates of the fixed effects and the conditional expectations of the random effects do not have the analytical solutions they do in the case of the linear model.

As shown in §2 parameter estimation for the Gaussian linear mixed-effects model can be re-expressed as a penalized least-squares problem. It is natural also to re-phrase Gaussian nonlinear mixed-effects models as penalized nonlinear least squares problems approached through linear approximation to the nonlinear model.

For a purely fixed-effects nonlinear regression model a common iterative parameter estimation method is the Gauss-Newton method (Bates and Watts, 1988, §2.2) wherein the nonlinear model  $\mathbf{f}(\boldsymbol{\beta})$  is replaced by a first-order Taylor series approximation about current estimates  $\boldsymbol{\beta}^{(k)}$  as

$$\mathbf{f}(\boldsymbol{\beta}) \approx \mathbf{f}(\boldsymbol{\beta}^{(k)}) + \left. \frac{d\mathbf{f}}{d\boldsymbol{\beta}'} \right|_{\boldsymbol{\beta}=\boldsymbol{\beta}^{(k)}} (\boldsymbol{\beta} - \boldsymbol{\beta}^{(k)})$$

The parameter increment  $\boldsymbol{\delta}^{(k)} = \boldsymbol{\beta}^{(k+1)} - \boldsymbol{\beta}^{(k)}$  for the  $k$ th iteration is calculated as the least squares solution of

$$\mathbf{y} - \mathbf{f}(\boldsymbol{\beta}^{(k)}) = \left. \frac{d\mathbf{f}}{d\boldsymbol{\beta}'} \right|_{\boldsymbol{\beta}=\boldsymbol{\beta}^{(k)}} (\boldsymbol{\beta} - \boldsymbol{\beta}^{(k)}) + \boldsymbol{\epsilon}^{(k)}$$

A similar iterative scheme can be used to determine the conditional estimates  $\hat{\boldsymbol{\beta}}(\boldsymbol{\theta})$  and the conditional modes of the distribution of the random effects given a value of  $\boldsymbol{\theta}$  in a nonlinear mixed effects models. Some care must be taken when incorporating the penalty terms such as  $\Delta \mathbf{b}_i$  into the nonlinear problem because the penalty term is a linear function of  $\mathbf{b}_i$ , not a linear function of the increment  $\mathbf{b}_i^{(k+1)} - \mathbf{b}_i^{(k)}$ . Two possible formulations of the penalized nonlinear least squares problem are

$$\begin{bmatrix} \mathbf{y}_i - \mathbf{f}_i(\boldsymbol{\beta}^{(k)}, \mathbf{b}_i^{(k)}) \\ -\Delta \mathbf{b}_i^{(k)} \end{bmatrix} = \begin{bmatrix} \mathbf{X}_i \\ \mathbf{0} \end{bmatrix} (\boldsymbol{\beta}^{(k+1)} - \boldsymbol{\beta}^{(k)}) + \begin{bmatrix} \mathbf{Z}_i \\ \Delta \end{bmatrix} (\mathbf{b}_i^{(k+1)} - \mathbf{b}_i^{(k)}) \quad i = 1, \dots, M$$

or

$$\begin{bmatrix} \mathbf{y}_i - \mathbf{f}_i(\boldsymbol{\beta}^{(k)}, \mathbf{b}_i^{(k)}) + \mathbf{X}_i \boldsymbol{\beta}^{(k)} + \mathbf{Z}_i \mathbf{b}_i^{(k)} \\ \mathbf{0} \end{bmatrix} = \begin{bmatrix} \mathbf{X}_i \\ \mathbf{0} \end{bmatrix} \boldsymbol{\beta}^{(k+1)} + \begin{bmatrix} \mathbf{Z}_i \\ \Delta \end{bmatrix} \mathbf{b}_i^{(k+1)} \quad (29)$$

where

$$\mathbf{X}_i = \frac{d\mathbf{f}}{d\boldsymbol{\beta}'} \quad \text{and} \quad \mathbf{Z}_i = \frac{d\mathbf{f}}{d\mathbf{b}_i'}$$

We prefer the second formulation (29) as the calculation of the increment more closely follows the form of the solution of the penalized least-squares problem in the linear case. This is essentially the form of the increment used in Lindstrom and Bates (1990) although this form allows generalization to a multiple levels of random effects.

Once updated values of the fixed-effects parameters  $\boldsymbol{\beta}$  and the conditional modes of the random effects are available, the derivative matrices  $\mathbf{X}_i$  and  $\mathbf{Z}_i$  are re-evaluated and the variance-covariance parameters  $\boldsymbol{\theta}$  are updated by several EM iterations. Only the EM iterations are used at this point because the  $\mathbf{X}_i$  and  $\mathbf{Z}_i$  matrices are going to be

recalculated the next time  $\beta$  and the conditional modes of the random effects are updated. It is not worthwhile determining the exact optimal  $\theta$  for the approximate problem that will subsequently be modified. Once the process of updating  $\beta$  and  $\theta$  separately stabilizes, they can be optimized jointly with higher order approximations to the log-likelihood for the nonlinear mixed-effects model as described in Pinheiro and Bates (1995). The Laplacian approximation would be a good choice of a higher-order approximation.

The fact that the profiled log-likelihood or profiled log-restricted-likelihood for a linear mixed-effects model can be quickly evaluated and also that EM iterations are very fast makes the early optimization of  $\beta$  and  $\theta$  for a nonlinear mixed-effects model faster and more stable. Also, the current formulation of the linear mixed-effects model allows extension of nonlinear mixed-effects models to multiple nested levels of random effects which was not previously available.

### 5.3 Implementation

The computational methods for multilevel Gaussian mixed-effects models, general linear mixed-effects models and nonlinear nested mixed-effects models described in this paper are implemented in version 3.0 of the NLME library for S, S-PLUS, and R. Documentation and source code for this library is available at

<http://franz.stat.wisc.edu/pub/NLME/>

## References

- Anderson, E., Bai, Z., Bischoff, C., Demmel, J., Dongarra, J., DuCroz, J., Greenbaum, A., Hammarling, S., McKenney, A., Ostrouchov, S. and Sorensen, D. (1994). *LA-PACK Users' Guide, 2nd ed.*, SIAM, Philadelphia.
- Bates, D. M. (1983). The derivative of  $|X'X|$  and its uses, *Technometrics* **25**: 373–376.
- Bates, D. M. and Watts, D. G. (1988). *Nonlinear Regression Analysis and Its Applications*, Wiley, New York.
- Davidian, M. and Giltinan, D. M. (1995). *Nonlinear Models for Repeated Measurement Data*, 1st edn, Chapman & Hall, London.
- Dongarra, J. J., Bunch, J. R., Moler, C. B. and Stewart, G. W. (1979). *Linpack Users' Guide*, SIAM, Philadelphia.
- Goldstein, H. (1995). *Multilevel Statistical Models*, Halstead Press, New York.
- Golub, G. H. and Pereyra, V. (1973). The differentiation of pseudo-inverses and nonlinear least squares problems whose variables separate, *J. SIAM* **10**: 413–432.
- Harville, D. A. (1976). Extension of the Gauss-Markov theorem to include the estimation of random effects, *The Annals of Statistics* **4**: 384–395.

- Kaufman, L. (1975). A variable projection method for solving separable nonlinear least squares problems, *BIT* **15**: 49–57.
- Laird, N. M. and Ware, J. H. (1982). Random-effects models for longitudinal data, *Biometrics* **38**: 963–974.
- Lindstrom, M. J. and Bates, D. M. (1988). Newton-Raphson and EM algorithms for linear mixed-effects models for repeated-measures data, *Journal of the American Statistical Association* **83**: 1014–1022.
- Lindstrom, M. J. and Bates, D. M. (1990). Nonlinear mixed effects models for repeated measures data, *Biometrics* **46**: 673–687.
- Liu, C. and Rubin, D. B. (1994). The ECME algorithm: A simple extension of EM and ECM with faster monotone convergence, *Biometrika* **81**: 633–648.
- Longford, N. T. (1993). *Random Coefficient Models*, Oxford University Press, New York.
- Pinheiro, J. C. and Bates, D. M. (1995). Approximations to the log-likelihood function in the nonlinear mixed-effects model, *Journal of Computational and Graphical Statistics* **4**(1): 12–35.
- Pinheiro, J. C. and Bates, D. M. (1996). Unconstrained parameterizations for variance-covariance matrices, *Statistics and Computing* **6**: 289–296.
- Thisted, R. A. (1988). *Elements of Statistical Computing*, Chapman & Hall, London.