

Penalized least squares versus generalized least squares representations of linear mixed models

Douglas Bates
Department of Statistics
University of Wisconsin – Madison

June 23, 2008

Abstract

The methods in the `lme4` package for `R` for fitting linear mixed models are based on sparse matrix methods, especially the Cholesky decomposition of sparse positive-semidefinite matrices, in a penalized least squares representation of the conditional model for the response given the random effects. The representation is similar to that in Henderson’s mixed-model equations. An alternative representation of the calculations is as a generalized least squares problem. We describe the two representations, show the equivalence of the two representations and explain why we feel that the penalized least squares approach is more versatile and more computationally efficient.

1 Definition of the model

We consider linear mixed models in which the random effects are represented by a q -dimensional random vector, \mathbf{B} , and the response is represented by an n -dimensional random vector, \mathbf{Y} . We observe a value, \mathbf{y} , of the response. The random effects are unobserved.

The marginal distribution of the random effects is a multivariate normal distribution with mean $\mathbf{0}$ and a variance-covariance matrix, $\mathbf{\Sigma}(\boldsymbol{\theta})$, that depends on a parameter vector, $\boldsymbol{\theta}$. Typically the dimension of $\boldsymbol{\theta}$ is much, much smaller than q . We will defer describing a particular parameterization until

later. For the time being we simply characterize the marginal distribution of \mathcal{B} as

$$\mathcal{B} \sim \mathcal{N}(\mathbf{0}, \Sigma(\theta)) \quad (1)$$

The conditional distribution, $\mathcal{Y}|\mathcal{B}$, is also multivariate normal. The conditional mean, $E[\mathcal{Y}|\mathcal{B} = \mathbf{b}]$, is a linear function of the p -dimensional fixed-effects parameter, β , and the q -dimensional random effects vector, \mathbf{b} , defined by the $n \times p$ and $n \times q$ model matrices \mathbf{X} and \mathbf{Z} as

$$E[\mathcal{Y}|\mathcal{B} = \mathbf{b}] = \mathbf{X}\beta + \mathbf{Z}\mathbf{b}. \quad (2)$$

The conditional variance-covariance matrix of \mathcal{Y} is simply $\sigma^2 \mathbf{I}_n$, where \mathbf{I}_n denotes the identity matrix of order n . Thus

$$\mathcal{Y}|\mathcal{B} \sim \mathcal{N}(\mathbf{X}\beta + \mathbf{Z}\mathbf{b}, \sigma^2 \mathbf{I}_n) \quad (3)$$

1.1 Variance-covariance of the random effects

The variance-covariance matrix, $\Sigma(\theta)$, of the random effects, \mathcal{B} , must be symmetric and positive semidefinite (i.e. $\mathbf{x}'\Sigma\mathbf{x} \geq 0, \forall \mathbf{x} \in \mathbb{R}^q$). Because the maximum likelihood estimate of a variance component can be zero, it is important to allow for a semidefinite Σ . That is, we do not assume that Σ is positive definite (i.e. $\mathbf{x}'\Sigma\mathbf{x} > 0, \forall \mathbf{x} \in \mathbb{R}^q$) and we do not assume that Σ^{-1} exists.

A positive semidefinite matrix such as Σ has a Cholesky decomposition of the so-called “LDL’” form. We use a slightly modified version

$$\Sigma(\theta) = \sigma^2 \mathbf{T}(\theta) \mathbf{S}(\theta) \mathbf{S}(\theta) \mathbf{T}(\theta)' \quad (4)$$

where σ is the same scale parameter that occurs in the variance-covariance of $\mathcal{Y}|\mathcal{B}$, $\mathbf{T}(\theta)$ is a unit lower-triangular $q \times q$ matrix and $\mathbf{S}(\theta)$ is a diagonal $q \times q$ matrix with nonnegative diagonal elements.

1.2 Orthogonal random effects

Let us define a q -dimensional random vector, \mathcal{U} , of orthogonal random effects with a marginal distribution of

$$\mathcal{U} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_q) \quad (5)$$

and express \mathbf{B} as a linear transformation of \mathbf{U} ,

$$\mathbf{B} = \mathbf{T}\mathbf{S}\mathbf{U}. \quad (6)$$

Note that the transformation (6) gives the desired distribution of \mathbf{B} in that $\mathbb{E}[\mathbf{B}] = \mathbf{T}\mathbf{S}\mathbb{E}[\mathbf{U}] = \mathbf{0}$ and

$$\text{Var}(\mathbf{B}) = \mathbb{E}[\mathbf{B}\mathbf{B}'] = \mathbf{T}\mathbf{S}\mathbb{E}[\mathbf{U}\mathbf{U}']\mathbf{S}\mathbf{T}' = \sigma^2\mathbf{T}\mathbf{S}\mathbf{S}\mathbf{T}' = \Sigma.$$

The conditional distribution, $\mathbf{Y}|\mathbf{U}$, can be derived from $\mathbf{Y}|\mathbf{B}$ as

$$\mathbf{Y}|\mathbf{U} \sim \mathcal{N}(\mathbf{X}\beta + \mathbf{Z}\mathbf{T}\mathbf{S}\mathbf{u}, \sigma^2\mathbf{I}) \quad (7)$$

We will write the transpose of $\mathbf{Z}\mathbf{T}\mathbf{S}$ as \mathbf{A} . Because the matrices \mathbf{T} and \mathbf{S} depend on the parameter $\boldsymbol{\theta}$, $\mathbf{A}(\boldsymbol{\theta})$ is also a function of $\boldsymbol{\theta}$. That is

$$\mathbf{A}'(\boldsymbol{\theta}) = \mathbf{Z}\mathbf{T}(\boldsymbol{\theta})\mathbf{S}(\boldsymbol{\theta}). \quad (8)$$

1.3 Sparse matrix methods

The reason for the peculiar definition of \mathbf{A} as the transpose of the model matrix is because \mathbf{A} is stored and manipulated as a sparse matrix. In the compressed column-oriented storage form for sparse matrices there are advantages to storing \mathbf{A} as a matrix of n columns and q rows. In particular, the CHOLMOD sparse matrix library allows us to evaluate the sparse Cholesky factor, \mathbf{L} , a lower triangular matrix that satisfies

$$\mathbf{L}(\boldsymbol{\theta})\mathbf{L}(\boldsymbol{\theta})' = \mathbf{P}(\mathbf{A}(\boldsymbol{\theta})\mathbf{A}(\boldsymbol{\theta})' + \mathbf{I}_q)\mathbf{P}', \quad (9)$$

directly from \mathbf{A} .

In (9) the $q \times q$ matrix \mathbf{P} is a “fill-reducing” permutation matrix determined from the pattern of nonzeros in the sparse model matrix \mathbf{Z} . It does not affect the statistical theory (if $\mathbf{U} \sim \mathcal{N}(\mathbf{0}, \sigma^2\mathbf{I})$ then $\mathbf{P}'\mathbf{U}$ also has a $\mathcal{N}(\mathbf{0}, \sigma^2\mathbf{I})$ distribution because $\mathbf{P}\mathbf{P}' = \mathbf{P}'\mathbf{P} = \mathbf{I}$) but, because it affects the number of nonzeros in \mathbf{L} , it can have a tremendous impact on the amount storage required for \mathbf{L} and the time required to evaluate it. Indeed, it is precisely because $\mathbf{L}(\boldsymbol{\theta})$ can be evaluated quickly, even for complex models applied the large data sets, that the `lmer` function is effective in fitting such models.

2 The penalized least squares approach to linear mixed models

Given a value of $\boldsymbol{\theta}$ we form $\mathbf{A}(\boldsymbol{\theta})$ from which we evaluate $\mathbf{L}(\boldsymbol{\theta})$. We can then solve for the $q \times p$ matrix, \mathbf{R}_{ZX} , in the system of equations

$$\mathbf{L}(\boldsymbol{\theta})\mathbf{R}_{ZX} = \mathbf{P}\mathbf{A}(\boldsymbol{\theta})\mathbf{X} \quad (10)$$

and for the $p \times p$ upper triangular matrix, \mathbf{R}_X , satisfying

$$\mathbf{R}_X'\mathbf{R}_X = \mathbf{X}'\mathbf{X} - \mathbf{R}_{ZX}'\mathbf{R}_{ZX} \quad (11)$$

The conditional mode, $\tilde{\mathbf{u}}(\boldsymbol{\theta})$, of the orthogonal random effects and the conditional mle, $\hat{\boldsymbol{\beta}}(\boldsymbol{\theta})$, of the fixed-effects parameters can be determined simultaneously as the solutions to a penalized least squares problem

$$\begin{bmatrix} \tilde{\mathbf{u}}(\boldsymbol{\theta}) \\ \hat{\boldsymbol{\beta}}(\boldsymbol{\theta}) \end{bmatrix} = \arg \min_{\mathbf{u}, \boldsymbol{\beta}} \left\| \begin{bmatrix} \mathbf{y} \\ \mathbf{0} \end{bmatrix} - \begin{bmatrix} \mathbf{A}'\mathbf{P}' & \mathbf{X} \\ \mathbf{I}_q & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{u} \\ \boldsymbol{\beta} \end{bmatrix} \right\|^2 \quad (12)$$

for which the solution satisfies

$$\begin{bmatrix} \mathbf{P}(\mathbf{A}\mathbf{A}' + \mathbf{I})\mathbf{P}' & \mathbf{P}\mathbf{A}\mathbf{X} \\ \mathbf{X}'\mathbf{A}'\mathbf{P}' & \mathbf{X}'\mathbf{X} \end{bmatrix} \begin{bmatrix} \tilde{\mathbf{u}}(\boldsymbol{\theta}) \\ \hat{\boldsymbol{\beta}}(\boldsymbol{\theta}) \end{bmatrix} = \begin{bmatrix} \mathbf{P}\mathbf{A}\mathbf{y} \\ \mathbf{X}'\mathbf{y} \end{bmatrix} \quad (13)$$

The Cholesky factor of the system matrix for the PLS problem can be expressed using \mathbf{L} , \mathbf{R}_{ZX} and \mathbf{R}_X because

$$\begin{bmatrix} \mathbf{P}(\mathbf{A}\mathbf{A}' + \mathbf{I})\mathbf{P}' & \mathbf{P}\mathbf{A}\mathbf{X} \\ \mathbf{X}'\mathbf{A}'\mathbf{P}' & \mathbf{X}'\mathbf{X} \end{bmatrix} = \begin{bmatrix} \mathbf{L} & \mathbf{0} \\ \mathbf{R}_{ZX}' & \mathbf{R}_X' \end{bmatrix} \begin{bmatrix} \mathbf{L}' & \mathbf{R}_{ZX} \\ \mathbf{0} & \mathbf{R}_X \end{bmatrix}. \quad (14)$$

In the `lme4` package the "mer" class is the representation of a mixed-effects model. The `A` slot contains the sparse matrix $\mathbf{A}(\boldsymbol{\theta})$ and the `L` slot contains the Cholesky factor $\mathbf{L}(\boldsymbol{\theta})$ satisfying (9). The `RZX` and `RX` slots contain $\mathbf{R}_{ZX}(\boldsymbol{\theta})$ and $\mathbf{R}_X(\boldsymbol{\theta})$ as dense matrices.

It is not necessary to solve for $\tilde{\mathbf{u}}(\boldsymbol{\theta})$ and $\hat{\boldsymbol{\beta}}(\boldsymbol{\theta})$ to evaluate the profiled log-likelihood as a function of $\boldsymbol{\theta}$, which is the log-likelihood evaluated at $\boldsymbol{\theta}$, $\hat{\boldsymbol{\beta}}(\boldsymbol{\theta})$ and $\hat{\sigma}^2(\boldsymbol{\theta})$. All that is needed for evaluation of the profiled log-likelihood is the penalized residual sum of squares, r^2 , and the determinant, $|\mathbf{A}\mathbf{A}' + \mathbf{I}| = |\mathbf{L}|^2$.

Because \mathbf{L} is triangular, its determinant is simply the product of its diagonal elements and, because $\mathbf{A}\mathbf{A}' + \mathbf{I}$ is positive definite, $|\mathbf{L}|^2 > 0$.

The profiled deviance (negative twice the profiled log-likelihood), as a function of $\boldsymbol{\theta}$ only ($\boldsymbol{\beta}$ and σ^2 at their conditional estimates), is

$$d(\boldsymbol{\theta}|\mathbf{y}) = \log(|\mathbf{L}|^2) + n \left(1 + \log(r^2) + \frac{2\pi}{n} \right) \quad (15)$$

The maximum likelihood estimates, $\hat{\boldsymbol{\theta}}$, satisfy

$$\hat{\boldsymbol{\theta}} = \arg \min_{\boldsymbol{\theta}} d(\boldsymbol{\theta}|\mathbf{y}) \quad (16)$$

Once the value of $\hat{\boldsymbol{\theta}}$ has been determined, the mle's of the other parameters are evaluated from (13) and

$$\hat{\sigma}^2(\boldsymbol{\theta}) = \frac{r^2}{n}. \quad (17)$$

2.1 Comments on the sparse matrix representation

Note that nothing has been said about the form of the sparse model matrix \mathbf{Z} other than the fact that it is sparse. The computational methods outlined above can be applied to models with multiple random effects terms in which the factors determining the random effects are nested or crossed or partially crossed.

3 The generalized least squares approach to linear mixed models

Another common approach to linear mixed models is to derive the marginal variance-covariance matrix of \mathbf{y} as a function of $\boldsymbol{\theta}$ and use that to determine the conditional estimates, $\hat{\boldsymbol{\beta}}(\boldsymbol{\theta})$, as the solution of a generalized least squares (GLS) problem. In the notation of §1 the marginal mean of \mathbf{y} is $E[\mathbf{y}] = \mathbf{X}\boldsymbol{\beta}$ and the marginal variance-covariance matrix is

$$\text{Var}(\mathbf{y}) = \sigma^2 (\mathbf{I}_n + \mathbf{ZTSST}'\mathbf{Z}') = \sigma^2 (\mathbf{I}_n + \mathbf{A}'\mathbf{A}) = \sigma^2 \mathbf{V}(\boldsymbol{\theta}), \quad (18)$$

where $\mathbf{V}(\boldsymbol{\theta}) = \mathbf{I}_n + \mathbf{A}'\mathbf{A}$.

The conditional estimates of β are often written as

$$\hat{\beta}(\theta) = (\mathbf{X}'\mathbf{V}^{-1}\mathbf{X})^{-1} \mathbf{X}'\mathbf{V}^{-1}\mathbf{y} \quad (19)$$

but, of course, this formula is not suitable for computation. The matrix $\mathbf{V}(\theta)$ is a symmetric $n \times n$ positive definite matrix and hence has a Cholesky factor. However, this factor is $n \times n$, not $q \times q$ and q is always smaller than n - sometimes orders of magnitude smaller.

3.1 Relating the GLS approach to the Cholesky factor L .

We can use the fact that

$$\mathbf{V}^{-1}(\theta) = (\mathbf{I}_n + \mathbf{A}'\mathbf{A})^{-1} = \mathbf{I}_n - \mathbf{A}'(\mathbf{I}_q + \mathbf{A}\mathbf{A}')^{-1}\mathbf{A} \quad (20)$$

to relate the GLS problem to the PLS problem. One way to establish (20) is simply to show that the product

$$\begin{aligned} (\mathbf{I} + \mathbf{A}'\mathbf{A}) \left(\mathbf{I} - \mathbf{A}'(\mathbf{I} + \mathbf{A}\mathbf{A}')^{-1}\mathbf{A} \right) \\ = \mathbf{I} + \mathbf{A}'\mathbf{A} - \mathbf{A}'(\mathbf{I} + \mathbf{A}\mathbf{A}')^{-1}(\mathbf{I} + \mathbf{A}\mathbf{A}')\mathbf{A} \\ = \mathbf{I} + \mathbf{A}'\mathbf{A} - \mathbf{A}'\mathbf{A} \\ = \mathbf{I}. \end{aligned}$$

Incorporating the permutation matrix \mathbf{P} we have

$$\begin{aligned} \mathbf{V}^{-1}(\theta) &= \mathbf{I}_n - \mathbf{A}'\mathbf{P}'\mathbf{P}(\mathbf{I}_q + \mathbf{A}\mathbf{A}')^{-1}\mathbf{P}'\mathbf{P}\mathbf{A} \\ &= \mathbf{I}_n - \mathbf{A}'\mathbf{P}'(\mathbf{L}\mathbf{L}')^{-1}\mathbf{P}\mathbf{A} \\ &= \mathbf{I}_n - (\mathbf{L}^{-1}\mathbf{P}\mathbf{A})' \mathbf{L}^{-1}\mathbf{P}\mathbf{A}. \end{aligned} \quad (21)$$

Even in this form we would not want to evaluate such a matrix but (21) does allow us to simplify many common expressions.

For example, the variance-covariance of the estimator $\hat{\beta}$, conditional on θ and σ , can be expressed as

$$\begin{aligned} \sigma^2 (\mathbf{X}'\mathbf{V}^{-1}(\theta)\mathbf{X})^{-1} &= \sigma^2 \left(\mathbf{X}'\mathbf{X} - (\mathbf{L}^{-1}\mathbf{P}\mathbf{A}\mathbf{X})' (\mathbf{L}^{-1}\mathbf{P}\mathbf{A}\mathbf{X}) \right)^{-1} \\ &= \sigma^2 (\mathbf{X}'\mathbf{X} - \mathbf{R}'_{ZX}\mathbf{R}_{ZX})^{-1} \\ &= \sigma^2 (\mathbf{R}'_X\mathbf{R}_X)^{-1}. \end{aligned} \quad (22)$$