

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

Douglas Bates  
Department of Statistics  
University of Wisconsin – Madison

April 6, 2022

## 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.

For our purposes, we will assume a “spherical” multivariate normal conditional distribution of  $\mathbf{Y}$ , given  $\mathbf{B}$ . That is, we assume the variance-covariance matrix of  $\mathbf{Y}|\mathbf{B}$  is simply  $\sigma^2\mathbf{I}_n$ , where  $\mathbf{I}_n$  denotes the identity matrix of order

$n$ . (The term “spherical” refers to the fact that contours of the conditional density are concentric spheres.)

The conditional mean,  $E[\mathbf{Y}|\mathcal{B} = \mathbf{b}]$ , is a linear function of  $\mathbf{b}$  and the  $p$ -dimensional fixed-effects parameter,  $\boldsymbol{\beta}$ ,

$$E[\mathbf{Y}|\mathcal{B} = \mathbf{b}] = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{b}, \quad (1)$$

where  $\mathbf{X}$  and  $\mathbf{Z}$  are known model matrices of sizes  $n \times p$  and  $n \times q$ , respectively. Thus

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

The marginal distribution of the random effects

$$\mathcal{B} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \boldsymbol{\Sigma}(\boldsymbol{\theta})) \quad (3)$$

is also multivariate normal, with mean  $\mathbf{0}$  and variance-covariance matrix  $\sigma^2 \boldsymbol{\Sigma}(\boldsymbol{\theta})$ . The scalar,  $\sigma^2$ , in (3) is the same as the  $\sigma^2$  in (2). As described in the next section, the relative variance-covariance matrix,  $\boldsymbol{\Sigma}(\boldsymbol{\theta})$ , is a  $q \times q$  positive semidefinite matrix depending on a parameter vector,  $\boldsymbol{\theta}$ . Typically the dimension of  $\boldsymbol{\theta}$  is much, much smaller than  $q$ .

## 1.1 Variance-covariance of the random effects

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

A positive semidefinite matrix such as  $\boldsymbol{\Sigma}$  has a Cholesky decomposition of the so-called “LDL” form. We use a slight modification of this form,

$$\boldsymbol{\Sigma}(\boldsymbol{\theta}) = \mathbf{T}(\boldsymbol{\theta})\mathbf{S}(\boldsymbol{\theta})\mathbf{S}(\boldsymbol{\theta})\mathbf{T}(\boldsymbol{\theta})', \quad (4)$$

where  $\mathbf{T}(\boldsymbol{\theta})$  is a unit lower-triangular  $q \times q$  matrix and  $\mathbf{S}(\boldsymbol{\theta})$  is a diagonal  $q \times q$  matrix with nonnegative diagonal elements that act as scale factors. (They are the relative standard deviations of certain linear combinations of the random effects.) Thus,  $\mathbf{T}$  is a triangular matrix and  $\mathbf{S}$  is a scale matrix.

Both  $\mathbf{T}$  and  $\mathbf{S}$  are highly patterned.

## 1.2 Orthogonal random effects

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

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

and, for a given value of  $\boldsymbol{\theta}$ , express  $\mathbf{B}$  as a linear transformation of  $\mathbf{U}$ ,

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

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

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

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

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

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  $\boldsymbol{\theta}$ ,  $\mathbf{A}$  is also a function of  $\boldsymbol{\theta}$ ,

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

In applications, the matrix  $\mathbf{Z}$  is derived from indicator columns of the levels of one or more factors in the data and is a *sparse* matrix, in the sense that most of its elements are zero. The matrix  $\mathbf{A}$  is also sparse. In fact, the structure of  $\mathbf{T}$  and  $\mathbf{S}$  are such that pattern of nonzeros in  $\mathbf{A}$  is that same as that in  $\mathbf{Z}'$ .

## 1.3 Sparse matrix methods

The reason for defining  $\mathbf{A}$  as the transpose of a model matrix is because  $\mathbf{A}$  is stored and manipulated as a sparse matrix. In the compressed column-oriented storage form that we use 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,  $\mathbf{L}(\boldsymbol{\theta})$ , a sparse 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  $\mathbf{A}(\boldsymbol{\theta})$ .

In (9) the  $q \times q$  matrix  $\mathbf{P}$  is a “fill-reducing” permutation matrix determined from the pattern of nonzeros in  $\mathbf{Z}$ .  $\mathbf{P}$  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  $\mathbf{L}$  from  $\mathbf{A}$ . 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. Several slots in this class are matrices corresponding directly to the matrices in the preceding equations. The `A` slot contains the sparse matrix  $\mathbf{A}(\boldsymbol{\theta})$  and the `L` slot contains the sparse Cholesky factor,  $\mathbf{L}(\boldsymbol{\theta})$ . The `RZX` and `RX` slots contain  $\mathbf{R}_{\mathbf{Z}\mathbf{X}}(\boldsymbol{\theta})$  and  $\mathbf{R}_{\mathbf{X}}(\boldsymbol{\theta})$ , respectively, stored as dense matrices.

It is not necessary to solve for  $\tilde{\boldsymbol{\mu}}(\boldsymbol{\theta})$  and  $\hat{\boldsymbol{\beta}}(\boldsymbol{\theta})$  to evaluate the *profiled* log-likelihood, which is the log-likelihood evaluated  $\boldsymbol{\theta}$  and the conditional estimates of the other parameters,  $\hat{\boldsymbol{\beta}}(\boldsymbol{\theta})$  and  $\hat{\sigma}^2(\boldsymbol{\theta})$ . All that is needed for evaluation of the profiled log-likelihood is the (penalized) residual sum of squares,  $r^2$ , from the penalized least squares problem (12) and the determinant  $|\mathbf{A}\mathbf{A}' + \mathbf{I}| = |\mathbf{L}|^2$ . Because  $\mathbf{L}$  is triangular, its determinant is easily evaluated as the product of its diagonal elements. Furthermore,  $|\mathbf{L}|^2 > 0$  because it is equal to  $|\mathbf{A}\mathbf{A}' + \mathbf{I}|$ , which is the determinant of a positive definite matrix. Thus  $\log(|\mathbf{L}|^2)$  is both well-defined and easily calculated from  $\mathbf{L}$ .

The profiled deviance (negative twice the profiled log-likelihood), as a function of  $\boldsymbol{\theta}$  only ( $\boldsymbol{\beta}$  and  $\sigma^2$  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 of  $\boldsymbol{\beta}$  is evaluated from (13) and the mle of  $\sigma^2$  as  $\hat{\sigma}^2(\boldsymbol{\theta}) = r^2/n$ .

Note that nothing has been said about the form of the sparse model matrix,  $\mathbf{Z}$ , other than the fact that it is sparse. In contrast to other methods for linear mixed models, these results apply to models where  $\mathbf{Z}$  is derived from crossed or partially crossed grouping factors, in addition to models with multiple, nested grouping factors.

The system (13) is similar to Henderson's "mixed-model equations" (reference?). One important difference between (13) and Henderson's formulation is that Henderson represented his system of equations in terms of  $\boldsymbol{\Sigma}^{-1}$  and, in important practical examples,  $\boldsymbol{\Sigma}^{-1}$  does not exist at the parameter estimates. Also, Henderson assumed that equations like (13) would need to be solved explicitly and, as we have seen, only the decomposition of the system matrix is needed for evaluation of the profiled log-likelihood. The same is

true of the profiled the logarithm of the REML criterion, which we define later.

### 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

$$\text{Var}(\mathbf{y}) = \sigma^2 (\mathbf{I}_n + \mathbf{Z}\mathbf{T}\mathbf{S}\mathbf{S}\mathbf{T}'\mathbf{Z}') = \sigma^2 (\mathbf{I}_n + \mathbf{A}'\mathbf{A}) = \sigma^2 \mathbf{V}(\boldsymbol{\theta}), \quad (17)$$

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

The conditional estimates of  $\boldsymbol{\beta}$  are often written as

$$\hat{\boldsymbol{\beta}}(\boldsymbol{\theta}) = (\mathbf{X}'\mathbf{V}^{-1}\mathbf{X})^{-1} \mathbf{X}'\mathbf{V}^{-1}\mathbf{y} \quad (18)$$

but, of course, this formula is not suitable for computation. The matrix  $\mathbf{V}(\boldsymbol{\theta})$  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  $n$  is always larger than  $q$  — sometimes orders of magnitude larger. Blithely writing a formula in terms of  $\mathbf{V}^{-1}$  when  $\mathbf{V}$  is  $n \times n$ , and  $n$  can be in the millions does not a computational formula make.

#### 3.1 Relating the GLS approach to the Cholesky factor

We can use the fact that

$$\mathbf{V}^{-1}(\boldsymbol{\theta}) = (\mathbf{I}_n + \mathbf{A}'\mathbf{A})^{-1} = \mathbf{I}_n - \mathbf{A}'(\mathbf{I}_q + \mathbf{A}\mathbf{A}')^{-1}\mathbf{A} \quad (19)$$

to relate the GLS problem to the PLS problem. One way to establish (19) is simply to show that the product

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

Incorporating the permutation matrix  $\mathbf{P}$  we have

$$\begin{aligned}
\mathbf{V}^{-1}(\boldsymbol{\theta}) &= \mathbf{I}_n - \mathbf{A}'\mathbf{P}'\mathbf{P}(\mathbf{I}_q + \mathbf{A}\mathbf{A}')^{-1}\mathbf{P}'\mathbf{P}\mathbf{A} \\
&= \mathbf{I}_n - \mathbf{A}'\mathbf{P}'(\mathbf{L}\mathbf{L}')^{-1}\mathbf{P}\mathbf{A} \\
&= \mathbf{I}_n - (\mathbf{L}^{-1}\mathbf{P}\mathbf{A})'\mathbf{L}^{-1}\mathbf{P}\mathbf{A}.
\end{aligned} \tag{20}$$

Even in this form we would not want to routinely evaluate  $\mathbf{V}^{-1}$ . However, (20) does allow us to simplify many common expressions.

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

$$\begin{aligned}
\sigma^2 (\mathbf{X}'\mathbf{V}^{-1}(\boldsymbol{\theta})\mathbf{X})^{-1} &= \sigma^2 \left( \mathbf{X}'\mathbf{X} - (\mathbf{L}^{-1}\mathbf{P}\mathbf{A}\mathbf{X})'(\mathbf{L}^{-1}\mathbf{P}\mathbf{A}\mathbf{X}) \right)^{-1} \\
&= \sigma^2 (\mathbf{X}'\mathbf{X} - \mathbf{R}'_{ZX}\mathbf{R}_{ZX})^{-1} \\
&= \sigma^2 (\mathbf{R}'_X\mathbf{R}_X)^{-1}.
\end{aligned} \tag{21}$$

## 4 Trace of the “hat” matrix

Another calculation that is of interest to some is the the trace of the “hat” matrix, which can be written as

$$\begin{aligned}
&\text{tr} \left( [\mathbf{A}' \quad \mathbf{X}] \left( \begin{bmatrix} \mathbf{A}' & \mathbf{X}' \\ \mathbf{I} & \mathbf{0} \end{bmatrix}' \begin{bmatrix} \mathbf{A}' & \mathbf{X}' \\ \mathbf{I} & \mathbf{0} \end{bmatrix} \right)^{-1} \begin{bmatrix} \mathbf{A} \\ \mathbf{X}' \end{bmatrix} \right) \\
&= \text{tr} \left( [\mathbf{A}' \quad \mathbf{X}] \left( \begin{bmatrix} \mathbf{L} & \mathbf{0} \\ \mathbf{R}'_{ZX} & \mathbf{R}'_X \end{bmatrix}' \begin{bmatrix} \mathbf{L}' & \mathbf{R}_{ZX} \\ \mathbf{0} & \mathbf{R}_X \end{bmatrix} \right)^{-1} \begin{bmatrix} \mathbf{A} \\ \mathbf{X}' \end{bmatrix} \right)
\end{aligned} \tag{22}$$