

Mixed models in R using the lme4 package

Part 6: Theory of linear mixed models, evaluating precision of estimates

Douglas Bates

University of Wisconsin - Madison
and R Development Core Team

[<Douglas.Bates@R-project.org>](mailto:Douglas.Bates@R-project.org)

University of Lausanne
July 2, 2009

Outline

Definition of linear mixed models

The penalized least squares problem

The sparse Cholesky factor

Evaluating the likelihood

Outline

Definition of linear mixed models

The penalized least squares problem

The sparse Cholesky factor

Evaluating the likelihood

Outline

Definition of linear mixed models

The penalized least squares problem

The sparse Cholesky factor

Evaluating the likelihood

Outline

Definition of linear mixed models

The penalized least squares problem

The sparse Cholesky factor

Evaluating the likelihood

Outline

Definition of linear mixed models

The penalized least squares problem

The sparse Cholesky factor

Evaluating the likelihood

Definition of linear mixed models

- As previously stated, we define a linear mixed model in terms of two random variables: the n -dimensional \mathbf{y} and the q -dimensional \mathbf{B}
- The probability model specifies the conditional distribution

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

and the unconditional distribution

$$\mathbf{B} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}(\boldsymbol{\theta})).$$

These distributions depend on the parameters $\boldsymbol{\beta}$, $\boldsymbol{\theta}$ and σ .

- The probability model defines the *likelihood* of the parameters, given the observed data, \mathbf{y} . In theory all we need to know is how to define the likelihood from the data so that we can maximize the likelihood with respect to the parameters. In practice we want to be able to evaluate it quickly and accurately.

Properties of $\Sigma(\theta)$; generating it

- Because it is a variance-covariance matrix, the $q \times q$ $\Sigma(\theta)$ must be symmetric and *positive semi-definite*, which means, in effect, that it has a “square root” — there must be another matrix that, when multiplied by its transpose, gives $\Sigma(\theta)$.
- We never really form Σ ; we always work with the *relative covariance factor*, $\Lambda(\theta)$, defined so that

$$\Sigma(\theta) = \sigma^2 \Lambda(\theta) \Lambda'(\theta)$$

where σ^2 is the same variance parameter as in $(\mathcal{Y}|\mathcal{B} = b)$.

- We also work with a q -dimensional “spherical” or “unit” random-effects vector, \mathcal{U} , such that

$$\mathcal{U} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_q), \quad \mathcal{B} = \Lambda(\theta) \mathcal{U} \Rightarrow \text{Var}(\mathcal{B}) = \sigma^2 \Lambda \Lambda' = \Sigma.$$

- The linear predictor expression becomes

$$\mathbf{Z}b + \mathbf{X}\beta = \mathbf{Z}\Lambda(\theta)\mathbf{u} + \mathbf{X}\beta = \mathbf{U}(\theta)\mathbf{u} + \mathbf{X}\beta$$

where $\mathbf{U}(\theta) = \mathbf{Z}\Lambda(\theta)$.

The conditional mean $\mu_{\mathbf{u}|\mathbf{y}=\mathbf{y}}$

- Although the probability model is defined from $(\mathbf{y}|\mathbf{u} = \mathbf{u})$, we observe \mathbf{y} , not \mathbf{u} (or \mathbf{b}) so we want to work with the other conditional distribution, $(\mathbf{u}|\mathbf{y} = \mathbf{y})$.
- The joint distribution of \mathbf{y} and \mathbf{u} is Gaussian with density

$$\begin{aligned} f_{\mathbf{y},\mathbf{u}}(\mathbf{y}, \mathbf{u}) &= f_{\mathbf{y}|\mathbf{u}}(\mathbf{y}|\mathbf{u}) f_{\mathbf{u}}(\mathbf{u}) \\ &= \frac{\exp(-\frac{1}{2}\|\mathbf{y} - \mathbf{X}\boldsymbol{\beta} - \mathbf{U}\mathbf{u}\|^2)}{(2\pi\sigma^2)^{n/2}} \frac{\exp(-\frac{1}{2}\|\mathbf{u}\|^2)}{(2\pi\sigma^2)^{q/2}} \\ &= \frac{\exp(-\frac{1}{2} [\|\mathbf{y} - \mathbf{X}\boldsymbol{\beta} - \mathbf{U}\mathbf{u}\|^2 + \|\mathbf{u}\|^2])}{(2\pi\sigma^2)^{(n+q)/2}} \end{aligned}$$

- $(\mathbf{u}|\mathbf{y} = \mathbf{y})$ is also Gaussian so its mean is its mode. I.e.

$$\mu_{\mathbf{u}|\mathbf{y}=\mathbf{y}} = \arg \min_{\mathbf{u}} \left[\|\mathbf{y} - \mathbf{X}\boldsymbol{\beta} - \mathbf{U}(\boldsymbol{\theta})\mathbf{u}\|^2 + \|\mathbf{u}\|^2 \right]$$

Outline

Definition of linear mixed models

The penalized least squares problem

The sparse Cholesky factor

Evaluating the likelihood

Minimizing a penalized sum of squared residuals

- An expression like $\|\mathbf{y} - \mathbf{X}\boldsymbol{\beta} - \mathbf{U}(\boldsymbol{\theta})\mathbf{u}\|^2 + \|\mathbf{u}\|^2$ is called a *penalized sum of squared residuals* because $\|\mathbf{y} - \mathbf{X}\boldsymbol{\beta} - \mathbf{U}(\boldsymbol{\theta})\mathbf{u}\|^2$ is a sum of squared residuals and $\|\mathbf{u}\|^2$ is a penalty on the size of the vector \mathbf{u} .
- Determining $\mu_{\mathbf{u}|\mathcal{Y}=\mathbf{y}}$ as the minimizer of this expression is a *penalized least squares* (PLS) problem. In this case it is a *penalized linear least squares problem* that we can solve directly (i.e. without iterating).
- One way to determine the solution is to rephrase it as a linear least squares problem for an extended residual vector

$$\mu_{\mathbf{u}|\mathcal{Y}=\mathbf{y}} = \arg \min_{\mathbf{u}} \left\| \begin{bmatrix} \mathbf{y} - \mathbf{X}\boldsymbol{\beta} \\ \mathbf{0} \end{bmatrix} - \begin{bmatrix} \mathbf{U}(\boldsymbol{\theta}) \\ \mathbf{I}_q \end{bmatrix} \mathbf{u} \right\|^2$$

This is sometimes called a *pseudo-data* approach because we create the effect of the penalty term, $\|\mathbf{u}\|^2$, by adding “pseudo-observations” to \mathbf{y} and to the predictor.

Solving the linear PLS problem

- The conditional mean satisfies the equations

$$[U(\boldsymbol{\theta})U'(\boldsymbol{\theta}) + \mathbf{I}_q]\boldsymbol{\mu}_{\mathcal{U}|\mathcal{Y}=\mathbf{y}} = U'(\mathbf{y} - \mathbf{X}\boldsymbol{\beta}).$$

- This would be interesting but not very important were it not for the fact that we actually can solve that system for $\boldsymbol{\mu}_{\mathcal{U}|\mathcal{Y}=\mathbf{y}}$ even when its dimension, q , is very, very large.
- Recall that $U(\boldsymbol{\theta}) = \mathbf{Z}\boldsymbol{\Lambda}(\boldsymbol{\theta})$. Because \mathbf{Z} is generated from indicator columns for the grouping factors, it is sparse. U is also very sparse.
- There are sophisticated and efficient ways of calculating a sparse Cholesky factor, which is a sparse, lower-triangular matrix $L(\boldsymbol{\theta})$ that satisfies

$$L(\boldsymbol{\theta})L'(\boldsymbol{\theta}) = U(\boldsymbol{\theta})'U(\boldsymbol{\theta}) + \mathbf{I}_q$$

and, from that, solving for $\boldsymbol{\mu}_{\mathcal{U}|\mathcal{Y}=\mathbf{y}}$.

Outline

Definition of linear mixed models

The penalized least squares problem

The sparse Cholesky factor

Evaluating the likelihood

The sparse Choleksy factor, $L(\theta)$

- Because the ability to evaluate the sparse Cholesky factor, $L(\theta)$, is the key to the computational methods in the `lme4` package, we consider this in detail.
- In practice we will evaluate $L(\theta)$ for many different values of θ when determining the ML or REML estimates of the parameters.
- As described in Davis (2006), §4.6, the calculation is performed in two steps: in the *symbolic decomposition* we determine the position of the nonzeros in L from those in U then, in the *numeric decomposition*, we determine the numerical values in those positions. Although the numeric decomposition may be done dozens, perhaps hundreds of times as we iterate on θ , the symbolic decomposition is only done once.

A fill-reducing permutation, P

- In practice it can be important while performing the symbolic decomposition to determine a *fill-reducing permutation*, which is written as a $q \times q$ permutation matrix, P . This matrix is just a re-ordering of the columns of I_q and has an orthogonality property, $PP' = P'P = I_q$.
- When P is used, the factor $L(\theta)$ is defined to be the sparse, lower-triangular matrix that satisfies

$$L(\theta)L'(\theta) = P [U'(\theta)U(\theta) + I_q] P'$$

- In the `Matrix` package for R , the `Cholesky` method for a sparse, symmetric matrix (class `dsCMatrix`) performs both the symbolic and numeric decomposition. By default, it determines a fill-reducing permutation, P . The `update` method for a Cholesky factor (class `CHMfactor`) performs the numeric decomposition only.

Applications to models with simple, scalar random effects

- Recall that, for a model with simple, scalar random-effects terms only, the matrix $\Sigma(\theta)$ is block-diagonal in k blocks and the i th block is $\sigma_i^2 \mathbf{I}_{n_i}$ where n_i is the number of levels in the i th grouping factor.
- The matrix $\Lambda(\theta)$ is also block-diagonal with the i th block being $\theta_i \mathbf{I}_{n_i}$, where $\theta_i = \sigma_i / \sigma$.
- Given the grouping factors for the model and a value of θ we produce \mathbf{U} then \mathbf{L} , using **Cholesky** the first time then **update**.
- To avoid recalculating we assign

flist a list of the grouping factors

nlev number of levels in each factor

Zt the transpose of the model matrix, \mathbf{Z}

theta current value of θ

Lambda current $\Lambda(\theta)$

Ut transpose of $\mathbf{U}(\theta) = \mathbf{Z}\Lambda(\theta)$

Cholesky factor for the Penicillin model

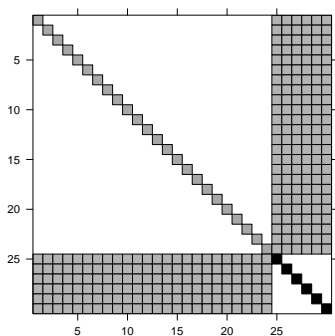
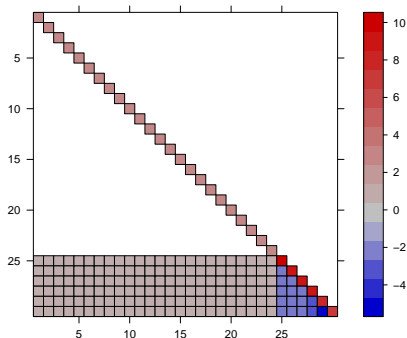
```
> flist <- subset(Penicillin, select = c(plate, sample))
> Zt <- do.call(rBind, lapply(flist, as, "sparseMatrix"))
> (nlev <- sapply(flist, function(f) length(levels(factor(f)))))

plate sample
    24      6

> theta <- c(1.2, 2.1)
> Lambda <- Diagonal(x = rep.int(theta, nlev))
> Ut <- crossprod(Lambda, Zt)
> str(L <- Cholesky(tcrossprod(Ut), LDL = FALSE, Imult = 1))
```

```
Formal class 'dCHMsimpl' [package "Matrix"] with 10 slots
 ..@ x      : num [1:189] 3.105 0.812 0.812 0.812 0.812 ...
 ..@ p      : int [1:31] 0 7 14 21 28 35 42 49 56 63 ...
 ..@ i      : int [1:189] 0 24 25 26 27 28 29 1 24 25 ...
 ..@ nz     : int [1:30] 7 7 7 7 7 7 7 7 7 7 ...
 ..@ nxt    : int [1:32] 1 2 3 4 5 6 7 8 9 10 ...
 ..@ prv    : int [1:32] 31 0 1 2 3 4 5 6 7 8 ...
 ..@ colcount: int [1:30] 7 7 7 7 7 7 7 7 7 7 ...
 ..@ perm   : int [1:30] 23 22 21 20 19 18 17 16 15 14 ...
 ..@ type   : int [1:4] 2 1 0 1
 ..@ Dim    : int [1:2] 30 30
```

Images of $U'U + I$ and L

 $U'U + I$  L

- Note that there are nonzeros in the lower right of L in positions that are zero in the lower triangle of $U'U + I$. This is described as “fill-in”.

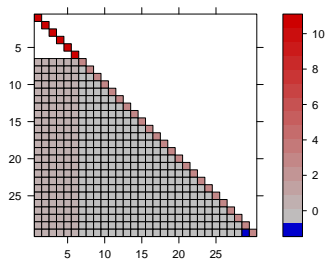
Reversing the order of the factors

- To show the effect of a fill-reducing permutation, we reverse the order of the factors and calculate the Cholesky factor with and without a fill-reducing permutation.
- We evaluate `nnzero` (number of nonzeros) for `L`, from the original factor order, and for `Lnoperm` and `Lperm`, the reversed factor order without and with permutation

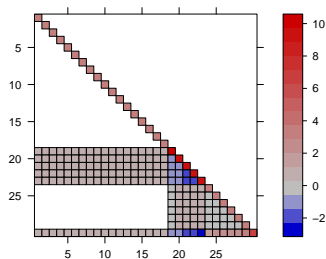
```
> Zt <- do.call(rBind, lapply(flist[2:1], as, "sparseMatrix"))
> Lambda <- Diagonal(x = rep.int(theta[2:1], nlev[2:1]))
> Ut <- crossprod(Lambda, Zt)
> Lnoperm <- Cholesky(tcrossprod(Ut), perm = FALSE, LDL = FALSE,
+   Imult = 1)
> Lperm <- Cholesky(tcrossprod(Ut), LDL = FALSE, Imult = 1)
> sapply(lapply(list(L, Lnoperm, Lperm), as, "sparseMatrix"),
+   nnzero)
```

```
[1] 189 450 204
```

Images of the reversed factor decompositions



L_{noperm}



L_{perm}

- Without permutation, we get the worse possible fill-in. With a fill-reducing permutation we get much less but still not as good as the original factor order.
- This is why the permutation is called “fill-reducing”, not “fill-minimizing”. Getting the fill-minimizing permutation in the general case is a very hard problem.

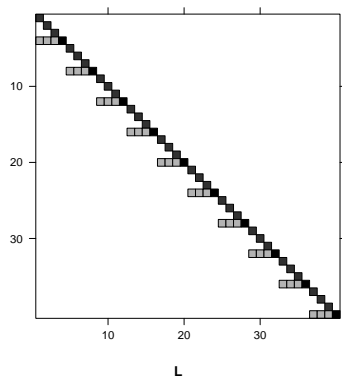
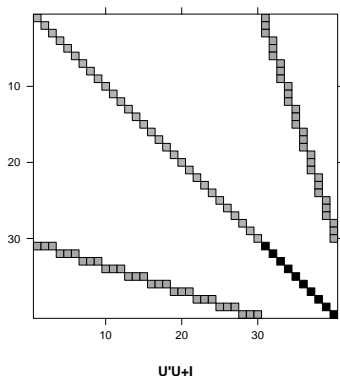
Cholesky factor for the Pastes data

- For the special case of nested grouping factors, such as in the `Pastes` and `classroom` data, there is no fill-in, regardless of the permutation.
- A permutation is nevertheless evaluated but it is a “post-ordering” that puts the nonzeros near the diagonal.

```
> Zt <- do.call(rBind, lapply(flist <- subset(Pastes,
+      , c(sample, batch)), as, "sparseMatrix"))
> nlev <- sapply(flist, function(f) length(levels(factor(f))))
> theta <- c(0.4, 0.5)
> Lambda <- Diagonal(x = rep.int(theta, nlev))
> Ut <- crossprod(Lambda, Zt)
> L <- Cholesky(tcrossprod(Ut), LDL = FALSE, Imult = 1)
> str(L@perm)
```

```
int [1:40] 2 1 0 30 5 4 3 31 8 7 ...
```

Image of the factor for the Pastes data



- The image for the Cholesky factor from the [classroom](#) data model is similar but, with more than 400 rows and columns, the squares for the nonzeros are difficult to see.

Outline

Definition of linear mixed models

The penalized least squares problem

The sparse Cholesky factor

Evaluating the likelihood

The conditional density, $f_{\mathbf{u}|\mathbf{y}=\mathbf{y}}$

- We know the joint density, $f_{\mathbf{y},\mathbf{u}}(\mathbf{y}, \mathbf{u})$, and

$$f_{\mathbf{u}|\mathbf{y}=\mathbf{y}}(\mathbf{u}|\mathbf{y}) = \frac{f_{\mathbf{y},\mathbf{u}}(\mathbf{y}, \mathbf{u})}{\int f_{\mathbf{y},\mathbf{u}}(\mathbf{y}, \mathbf{u}) d\mathbf{u}}$$

so we almost have $f_{\mathbf{u}|\mathbf{y}=\mathbf{y}}$. The trick is evaluating the integral in the denominator, which, it turns out, is exactly the likelihood, $L(\boldsymbol{\theta}, \boldsymbol{\beta}, \sigma^2|\mathbf{y})$, that we want to maximize.

- The Cholesky factor, $\mathbf{L}(\boldsymbol{\theta})$ is the key to doing this because

$$\mathbf{P}'\mathbf{L}(\boldsymbol{\theta})\mathbf{L}'(\boldsymbol{\theta})\mathbf{P}\boldsymbol{\mu}_{\mathbf{u}|\mathbf{y}=\mathbf{y}} = \mathbf{U}'(\mathbf{y} - \mathbf{X}\boldsymbol{\beta}).$$

Although the `Matrix` package provides a one-step `solve` method for this, we write it in stages:

Solve $\mathbf{L}\mathbf{c}_{\mathbf{u}} = \mathbf{P}\mathbf{U}'(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})$ for $\mathbf{c}_{\mathbf{u}}$.

Solve $\mathbf{L}'\mathbf{P}\boldsymbol{\mu} = \mathbf{c}_{\mathbf{u}}$ for $\mathbf{P}\boldsymbol{\mu}$ and $\boldsymbol{\mu}$ as $\mathbf{P}'\mathbf{P}\boldsymbol{\mu}$.

Define $r^2(\boldsymbol{\theta}, \boldsymbol{\beta}) = \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta} - \mathbf{U}\boldsymbol{\mu}_{\mathbf{u}|\mathbf{y}}\|^2 + \|\boldsymbol{\mu}_{\mathbf{u}|\mathbf{y}}\|^2$. Then

$$\|\mathbf{y} - \mathbf{X}\boldsymbol{\beta} - \mathbf{U}\mathbf{u}\|^2 + \|\mathbf{u}\|^2 = r^2(\boldsymbol{\theta}, \boldsymbol{\beta}) + \|\mathbf{c}_{\mathbf{u}} - \mathbf{L}'\mathbf{P}\mathbf{u}\|^2$$

Penalized sum of squared residuals



Recap

- For a linear mixed model, even one with a huge number of observations and random effects like the model for the grade point scores, evaluation of the ML or REML profiled deviance, given a value of θ , is straightforward. It involves updating T and S , then updating A , L , R_{ZX} , R_X , calculating the penalized residual sum of squares, r and a couple of determinants of triangular matrices.
- The profiled deviance can be optimized as a function of θ only. The dimension of θ is usually very small. For the grade point scores there are only three components to θ .