

# Mixed models in R using the lme4 package

## Part 4: Theory of linear mixed models

Douglas Bates

Madison

January 11, 2011

## Contents

### 1 Definition of linear mixed models

#### Definition of linear mixed models

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

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

and the unconditional distribution

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

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

- 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 $\boldsymbol{\Sigma}_\theta$ ; generating it

- Because it is a variance-covariance matrix, the  $q \times q$   $\boldsymbol{\Sigma}_\theta$  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  $\boldsymbol{\Sigma}_\theta$ .
- We never really form  $\boldsymbol{\Sigma}$ ; we always work with the *relative covariance factor*,  $\boldsymbol{\Lambda}_\theta$ , defined so that

$$\boldsymbol{\Sigma}_\theta = \sigma^2 \boldsymbol{\Lambda}_\theta \boldsymbol{\Lambda}_\theta^\top$$

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

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

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

- The linear predictor expression becomes

$$\mathbf{Zb} + \mathbf{X}\beta = \mathbf{Z}\mathbf{\Lambda}_\theta \mathbf{u} + \mathbf{X}\beta$$

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

- 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

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

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

$$\mu_{\mathbf{u}|\mathbf{y}} = \arg \min_{\mathbf{u}} \left[ \|\mathbf{y} - \mathbf{X}\beta - \mathbf{Z}\mathbf{\Lambda}_\theta \mathbf{u}\|^2 + \|\mathbf{u}\|^2 \right]$$

## 2 The penalized least squares problem

### Minimizing a penalized sum of squared residuals

- An expression like  $\|\mathbf{y} - \mathbf{X}\beta - \mathbf{Z}\mathbf{\Lambda}_\theta \mathbf{u}\|^2 + \|\mathbf{u}\|^2$  is called a *penalized sum of squared residuals* because  $\|\mathbf{y} - \mathbf{X}\beta - \mathbf{Z}\mathbf{\Lambda}_\theta \mathbf{u}\|^2$  is a sum of squared residuals and  $\|\mathbf{u}\|^2$  is a penalty on the size of the vector  $\mathbf{u}$ .
- Determining  $\mu_{\mathbf{u}|\mathbf{y}}$  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

$$\mu_{\mathbf{u}|\mathbf{y}} = \arg \min_{\mathbf{u}} \left\| \begin{bmatrix} \mathbf{y} - \mathbf{X}\beta \\ \mathbf{0} \end{bmatrix} - \begin{bmatrix} \mathbf{Z}\mathbf{\Lambda}_\theta \\ \mathbf{I}_q \end{bmatrix} \mathbf{u} \right\|^2$$

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

$$\left(\mathbf{Z}\Lambda_\theta\Lambda_\theta^\top\mathbf{Z}^\top + \mathbf{I}_q\right)\boldsymbol{\mu}_{\mathcal{U}|\mathcal{Y}} = \Lambda_\theta^\top\mathbf{Z}^\top(\mathbf{y} - \mathbf{X}\boldsymbol{\beta}).$$

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

$$\mathbf{L}_\theta\mathbf{L}_\theta^\top = \Lambda_\theta^\top\mathbf{Z}^\top\mathbf{Z}\Lambda_\theta + \mathbf{I}_q$$

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

## 3 The sparse Cholesky factor

### The sparse Cholesky factor, $\mathbf{L}_\theta$

- Because the ability to evaluate the sparse Cholesky factor,  $\mathbf{L}_\theta$ , is the key to the computational methods in the `lme4` package, we consider this in detail.
- In practice we will evaluate  $\mathbf{L}_\theta$  for many different values of  $\boldsymbol{\theta}$  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  $\mathbf{L}$  from those in  $\mathbf{Z}\Lambda_\theta$  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  $\boldsymbol{\theta}$ , the symbolic decomposition is only done once.

### A fill-reducing permutation, $\mathbf{P}$

- 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,  $\mathbf{P}$ . This matrix is just a re-ordering of the columns of  $\mathbf{I}_q$  and has an orthogonality property,  $\mathbf{P}\mathbf{P}^\top = \mathbf{P}^\top\mathbf{P} = \mathbf{I}_q$ .
- When  $\mathbf{P}$  is used, the factor  $\mathbf{L}_\theta$  is defined to be the sparse, lower-triangular matrix that satisfies

$$\mathbf{L}_\theta\mathbf{L}_\theta^\top = \mathbf{P} \left[ \Lambda_\theta^\top\mathbf{Z}^\top\mathbf{Z}\Lambda_\theta + \mathbf{I}_q \right] \mathbf{P}^\top$$

- 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,  $\mathbf{P}$ . 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  $\theta$  we produce  $\mathbf{Z}\Lambda_\theta$  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  $\theta$

`Lambda` current  $\Lambda_\theta$

`Ut` transpose of  $\mathbf{Z}\Lambda_\theta$

## Cholesky factor for the Penicillin model

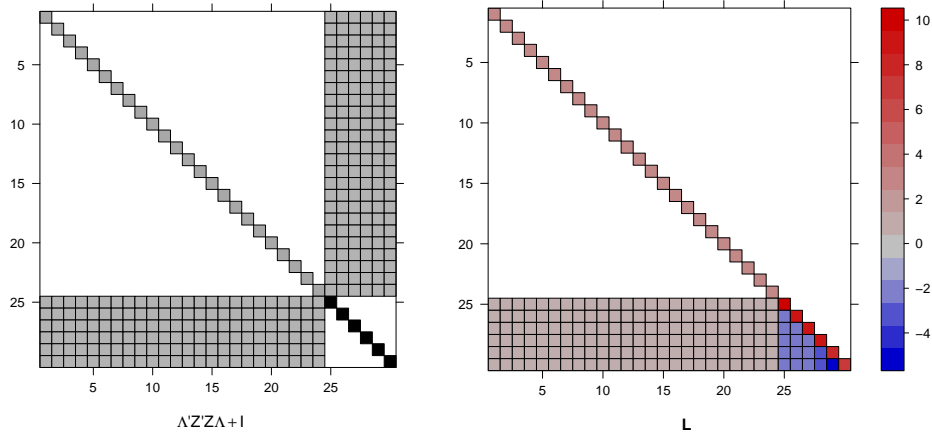
```
> flist <- subset(Penicillin, select = c(plate, sample))
> Zt <- do.call(rBind, lapply(flist, as, "sparseMatrix"))
> (nlev <- sapply(flist, function(f) length(levels(factor(f)))))
```

```
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 $\Lambda_\theta^\top Z^\top Z \Lambda_\theta + I$ and $L$



- Note that there are nonzeros in the lower right of  $L$  in positions that are zero in the lower triangle of  $\Lambda_\theta^\top Z^\top Z \Lambda_\theta + I_q$ . This is described as “fill-in”.

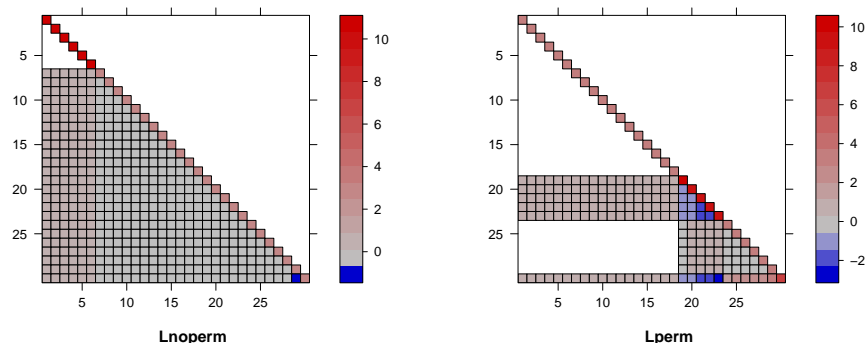
## Reversing the order of the factors

- 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  $L_{\text{noperm}}$  and  $L_{\text{perm}}$ , 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



- Without permutation, we get the worst 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.

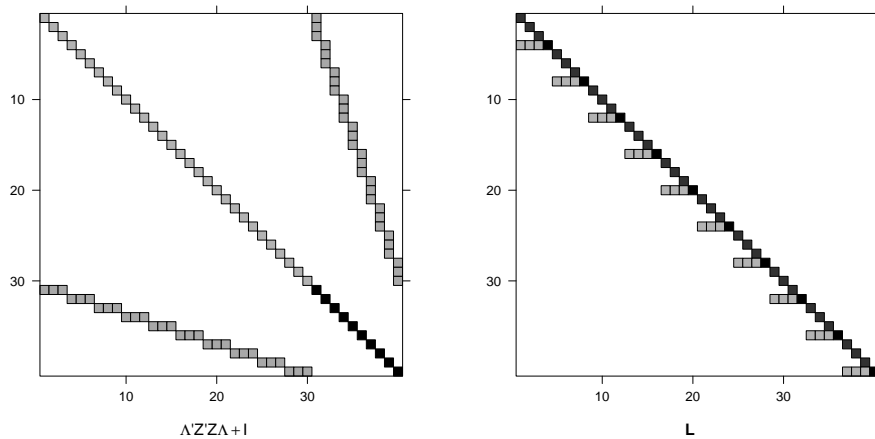
### Cholesky factor for the Pastes data

- 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.

## 4 Evaluating the likelihood

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

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

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

so we almost have  $f_{\mathbf{u}|\mathbf{y}}$ . 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,  $\mathbf{L}_\theta$  is the key to doing this because

$$\mathbf{P}^\top \mathbf{L}_\theta \mathbf{L}_\theta^\top \mathbf{P} \boldsymbol{\mu}_{\mathbf{u}|\mathbf{y}} = \boldsymbol{\Lambda}_\theta^\top \mathbf{Z}^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}).$$

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

1. Solve  $\mathbf{L}\mathbf{c}_u = \mathbf{P}\boldsymbol{\Lambda}_\theta^\top \mathbf{Z}^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})$  for  $\mathbf{c}_u$ .
2. Solve  $\mathbf{L}^\top \mathbf{P}\boldsymbol{\mu} = \mathbf{c}_u$  for  $\mathbf{P}\boldsymbol{\mu}_{\mathbf{u}|\mathbf{y}}$  and  $\boldsymbol{\mu}_{\mathbf{u}|\mathbf{y}}$  as  $\mathbf{P}^\top \mathbf{P}\boldsymbol{\mu}_{\mathbf{u}|\mathbf{y}}$ .

### Evaluating the likelihood

- The exponent of  $f_{\mathbf{y}, \mathbf{u}}(\mathbf{y}, \mathbf{u})$  can now be written

$$\|\mathbf{y} - \mathbf{X}\boldsymbol{\beta} - \mathbf{Z}\boldsymbol{\Lambda}_\theta \mathbf{u}\|^2 + \|\mathbf{u}\|^2 = r^2(\boldsymbol{\theta}, \boldsymbol{\beta}) + \|\mathbf{L}^\top \mathbf{P}(\mathbf{u} - \boldsymbol{\mu}_{\mathbf{u}|\mathbf{y}})\|^2.$$

where  $r^2(\boldsymbol{\theta}, \boldsymbol{\beta}) = \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta} - \mathbf{U}\boldsymbol{\mu}_{\mathbf{u}|\mathbf{y}}\|^2 + \|\boldsymbol{\mu}_{\mathbf{u}|\mathbf{y}}\|^2$ . The first term doesn't depend on  $\mathbf{u}$  and the second is relatively easy to integrate.

- Use the change of variable  $\mathbf{v} = \mathbf{L}^\top \mathbf{P}(\mathbf{u} - \boldsymbol{\mu}_{\mathbf{u}|\mathbf{y}})$ , with  $d\mathbf{v} = \text{abs}(|\mathbf{L}||\mathbf{P}|) d\mathbf{u}$ , in

$$\begin{aligned} \int \frac{\exp\left(\frac{-\|\mathbf{L}^\top \mathbf{P}(\mathbf{u} - \boldsymbol{\mu}_{\mathbf{u}|\mathbf{y}})\|^2}{2\sigma^2}\right)}{(2\pi\sigma^2)^{q/2}} d\mathbf{u} \\ = \int \frac{\exp\left(\frac{-\|\mathbf{v}\|^2}{2\sigma^2}\right)}{(2\pi\sigma^2)^{q/2}} \frac{d\mathbf{v}}{\text{abs}(|\mathbf{L}||\mathbf{P}|)} = \frac{1}{\text{abs}(|\mathbf{L}||\mathbf{P}|)} = \frac{1}{|\mathbf{L}|} \end{aligned}$$

because  $\text{abs}|\mathbf{P}| = 1$  and  $\text{abs}|\mathbf{L}|$ , which is the product of its diagonal elements, all of which are positive, is positive.

### Evaluating the likelihood (cont'd)

- As is often the case, it is easiest to write the log-likelihood. On the deviance scale (negative twice the log-likelihood)  $\ell(\boldsymbol{\theta}, \boldsymbol{\beta}, \sigma|\mathbf{y}) = \log L(\boldsymbol{\theta}, \boldsymbol{\beta}, \sigma|\mathbf{y})$  becomes

$$-2\ell(\boldsymbol{\theta}, \boldsymbol{\beta}, \sigma|\mathbf{y}) = n \log(2\pi\sigma^2) + \frac{r^2(\boldsymbol{\theta}, \boldsymbol{\beta})}{\sigma^2} + \log(|\mathbf{L}_\theta|^2)$$

- We wish to minimize the deviance. Its dependence on  $\sigma$  is straightforward. Given values of the other parameters, we can evaluate the conditional estimate

$$\widehat{\sigma^2}(\boldsymbol{\theta}, \boldsymbol{\beta}) = \frac{r^2(\boldsymbol{\theta}, \boldsymbol{\beta})}{n}$$

producing the *profiled deviance*

$$-2\tilde{\ell}(\boldsymbol{\theta}, \boldsymbol{\beta}|\mathbf{y}) = \log(|\mathbf{L}_\theta|^2) + n \left[ 1 + \log\left(\frac{2\pi r^2(\boldsymbol{\theta}, \boldsymbol{\beta})}{n}\right) \right]$$

- However, an even greater simplification is possible because the deviance depends on  $\beta$  only through  $r^2(\theta, \beta)$ .

### Profiling the deviance with respect to $\beta$

- Because the deviance depends on  $\beta$  only through  $r^2(\theta, \beta)$  we can obtain the conditional estimate,  $\hat{\beta}_\theta$ , by extending the PLS problem to

$$r_\theta^2 = \min_{\mathbf{u}, \beta} \left[ \|\mathbf{y} - \mathbf{X}\beta - \mathbf{Z}\Lambda_\theta \mathbf{u}\|^2 + \|\mathbf{u}\|^2 \right]$$

with the solution satisfying the equations

$$\begin{bmatrix} \Lambda_\theta^\top \mathbf{Z}^\top \mathbf{Z} \Lambda_\theta + \mathbf{I}_q & \mathbf{U}_\theta^\top \mathbf{X} \\ \mathbf{X}^\top \mathbf{Z} \Lambda_\theta & \mathbf{X}^\top \mathbf{X} \end{bmatrix} \begin{bmatrix} \mu_{\mathbf{u}|\mathbf{y}} \\ \hat{\beta}_\theta \end{bmatrix} = \begin{bmatrix} \Lambda_\theta^\top \mathbf{Z}^\top \mathbf{y} \\ \mathbf{X}^\top \mathbf{y} \end{bmatrix}$$

- The profiled deviance, which is a function of  $\theta$  only, is

$$-2\tilde{\ell}(\theta) = \log(|\mathbf{L}_\theta|^2) + n \left[ 1 + \log \left( \frac{2\pi r_\theta^2}{n} \right) \right]$$

### Solving the extended PLS problem

- For brevity we will no longer show the dependence of matrices and vectors on the parameter  $\theta$ .
- As before we use the sparse Cholesky decomposition, with  $\mathbf{L}$  and  $\mathbf{P}$  satisfying  $\mathbf{L}\mathbf{L}^\top = \mathbf{P}(\Lambda_\theta^\top \mathbf{Z}^\top \mathbf{Z} \Lambda_\theta + \mathbf{I})$  and  $\mathbf{c}_u$ , the solution to  $\mathbf{L}\mathbf{c}_u = \mathbf{P}\Lambda_\theta^\top \mathbf{Z}^\top \mathbf{y}$ .
- We extend the decomposition with the  $q \times p$  matrix  $\mathbf{R}_{ZX}$ , the upper triangular  $p \times p$  matrix  $\mathbf{R}_X$ , and the  $p$ -vector  $\mathbf{c}_\beta$  satisfying

$$\begin{aligned} \mathbf{L}\mathbf{R}_{ZX} &= \mathbf{P}\Lambda_\theta^\top \mathbf{Z}^\top \mathbf{X} \\ \mathbf{R}_X^\top \mathbf{R}_X &= \mathbf{X}^\top \mathbf{X} - \mathbf{R}_{ZX}^\top \mathbf{R}_{ZX} \\ \mathbf{R}_X^\top \mathbf{c}_\beta &= \mathbf{X}^\top \mathbf{y} - \mathbf{R}_{ZX}^\top \mathbf{c}_u \end{aligned}$$

so that

$$\begin{bmatrix} \mathbf{P}^\top \mathbf{L} & \mathbf{0} \\ \mathbf{R}_{ZX}^\top & \mathbf{R}_X^\top \end{bmatrix} \begin{bmatrix} \mathbf{L}^\top \mathbf{P} & \mathbf{R}_{ZX} \\ \mathbf{0} & \mathbf{R}_X \end{bmatrix} = \begin{bmatrix} \Lambda_\theta^\top \mathbf{Z}^\top \mathbf{Z} \Lambda_\theta + \mathbf{I} & \Lambda_\theta^\top \mathbf{Z}^\top \mathbf{X} \\ \mathbf{X}^\top \mathbf{Z} \Lambda_\theta & \mathbf{X}^\top \mathbf{X} \end{bmatrix}.$$

### Solving the extended PLS problem (cont'd)

- Finally we solve

$$\begin{aligned} \mathbf{R}_X \hat{\beta}_\theta &= \mathbf{c}_\beta \\ \mathbf{L}^\top \mathbf{P} \mu_{\mathbf{u}|\mathbf{y}} &= \mathbf{c}_u - \mathbf{R}_{ZX} \hat{\beta}_\theta \end{aligned}$$



- The profiled REML criterion also can be expressed simply. The criterion is

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

The same change-of-variable technique for evaluating the integral w.r.t.  $\mathbf{u}$  as  $1/\text{abs}(|\mathbf{L}|)$  produces  $1/\text{abs}(|\mathbf{R}_X|)$  here and removes  $(2\pi\sigma^2)^{p/2}$  from the denominator. On the deviance scale, the profiled REML criterion is

$$-2\tilde{\ell}_R(\boldsymbol{\theta}) = \log(|\mathbf{L}|^2) + \log(|\mathbf{R}_X|^2) + (n-p) \left[ 1 + \log\left(\frac{2\pi r_{\theta}^2}{n-p}\right) \right]$$

- These calculations can be expressed in a few lines of *R* code. Assume `rho` contains  $\mathbf{y}$ ,  $\mathbf{X}$ ,  $\mathbf{Z}_t$ , REML,  $\mathbf{L}$ , `nlev` and `XtX` ( $\mathbf{X}^T \mathbf{X}$ ).

### Code for evaluating the profiled deviance

```
profDev <- function(rho, theta) {
  stopifnot(is.numeric(theta), length(theta)==length(rho$nlev))
  Ut <- crossprod(Diagonal(x=rep.int(theta,rho$nlev)),rho$Zt)
  L <- update(rho$L, Ut, mult = 1)
  cu <- solve(L, solve(L, Ut %*% rho$y, sys = "P"), sys = "L")
  RZX <- solve(L, solve(L, Ut %*% rho$X, sys = "P"), sys = "L")
  RX <- chol(rho$XtX - crossprod(RZX))
  cb <- solve(t(RX),crossprod(rho$X,rho$y)- crossprod(RZX, cu))
  beta <- solve(RX, cb)
  u <- solve(L,solve(L,cu - RZX %*% beta, sys="Lt"), sys="Pt")
  fitted <- as.vector(crossprod(Ut, u) + rho$X %*% beta)
  prss <- sum(c(rho$y - fitted, as.vector(u))^2)
  n <- length(fitted); p <- ncol(RX)
  if (rho$REML) return(determinant(L)$mod +
    2 * determinant(RX)$mod +
    (n-p) * (1+log(2*pi*prss/(n-p))))
  determinant(L)$mod + n * (1 + log(2*pi*prss/n))
}
```

### Checking profDev, lmer version of fit

```
> invisible(lmer(mathgain ~ mathkind + minority + ses + (1|classid) + (1|schoolid), classroom,
```

```
npt = 5 , n = 2
rhobeg = 0.2 , rhoend = 2e-07
  0.020: 11:      11396.1;0.465012 0.298249
  0.0020: 17:      11391.5;0.340940 0.314970
  0.00020: 19:      11391.5;0.336337 0.313884
  2.0e-05: 23:      11391.5;0.336012 0.314470
  2.0e-06: 26:      11391.5;0.335979 0.314536
  2.0e-07: 29:      11391.5;0.335979 0.314536
At return
 34:      11391.532: 0.335980 0.314536
```

## How lmer works

- Essentially `lmer` takes its arguments and creates a structure like the `rho` environment shown above. The optimization of the profiled deviance or the profiled REML criterion happens within this environment.
- The creation of  $\mathbf{\Lambda}_\theta$  is somewhat more complex for models with vector-valued random effects but not excessively so.
- Some care is taken to avoid allocating storage for large objects during each function evaluation. Many of the objects created in `profDev` are updated in place within `lmer`.
- Once the optimizer, `bobyqa`, has converged some additional information for the summary is calculated.

## Summary

- 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  $\theta$ , is straightforward. It involves updating  $\mathbf{\Lambda}_\theta$ ,  $\mathbf{L}_\theta$ ,  $\mathbf{R}_{ZX}$ ,  $\mathbf{R}_X$ , calculating the penalized residual sum of squares,  $r_\theta^2$  and two determinants of triangular matrices.
- The profiled deviance can be optimized as a function of  $\theta$  only. The dimension of  $\theta$  is usually very small. For the grade point scores there are only three components to  $\theta$ .