

Mixed models in R using the lme4 package

Part 4: Theory of linear mixed models

Douglas Bates

8th International Amsterdam Conference
on Multilevel Analysis
<Bates@R-project.org>

2011-03-16

1 Definition of linear mixed models

Outline

- 1 Definition of linear mixed models
- 2 The penalized least squares problem

Outline

- 1 Definition of linear mixed models
- 2 The penalized least squares problem
- 3 The sparse Cholesky factor

Outline

- 1 Definition of linear mixed models
- 2 The penalized least squares problem
- 3 The sparse Cholesky factor
- 4 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

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

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 Σ_θ ; generating it

- Because it is a variance-covariance matrix, the $q \times q$ Σ_θ 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 Σ_θ .
- We never really form Σ ; we always work with the *relative covariance factor*, Λ_θ , defined so that

$$\Sigma_\theta = \sigma^2 \Lambda_\theta \Lambda_\theta^\top$$

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

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

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

- The linear predictor expression becomes

$$\mathbf{Z}\mathbf{b} + \mathbf{X}\boldsymbol{\beta} = \mathbf{Z}\Lambda_\theta \mathbf{u} + \mathbf{X}\boldsymbol{\beta}$$

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

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

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

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

Minimizing a penalized sum of squared residuals

- An expression like $\|\mathbf{y} - \mathbf{X}\boldsymbol{\beta} - \mathbf{Z}\boldsymbol{\Lambda}_\theta\mathbf{u}\|^2 + \|\mathbf{u}\|^2$ is called a *penalized sum of squared residuals* because $\|\mathbf{y} - \mathbf{X}\boldsymbol{\beta} - \mathbf{Z}\boldsymbol{\Lambda}_\theta\mathbf{u}\|^2$ 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

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

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

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

- This would be interesting but not very important were it not for the fact that we actually can solve that system for $\mu_{\mathcal{U}|\mathcal{Y}}$ 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}_θ 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 $\mu_{\mathcal{U}|\mathcal{Y}}$.

The sparse Choleksy factor, \mathbf{L}_θ

- Because the ability to evaluate the sparse Cholesky factor, \mathbf{L}_θ , is the key to the computational methods in the `lme4` package, we consider this in detail.
- In practice we will evaluate \mathbf{L}_θ 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 \mathbf{L} from those in $\mathbf{Z}\mathbf{\Lambda}_\theta$ 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 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, $PP^T = P^T P = I_q$.
- When P is used, the factor L_θ is defined to be the sparse, lower-triangular matrix that satisfies

$$L_\theta L_\theta^T = P [\Lambda_\theta^T Z_\theta^T Z_\theta \Lambda_\theta + I_q] P^T$$

- 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.

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}_θ 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:

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 .

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 σ 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 $\boldsymbol{\beta}$ only through $r^2(\boldsymbol{\theta}, \boldsymbol{\beta})$.

Profiling the deviance with respect to β

- Because the deviance depends on β 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 θ 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 θ .
- As before we use the sparse Cholesky decomposition, with L and P satisfying $LL^T = P(\Lambda_\theta^T Z^T Z \Lambda_\theta + I)$ and c_u , the solution to $Lc_u = P\Lambda_\theta^T Z^T y$.
- We extend the decomposition with the $q \times p$ matrix R_{ZX} , the upper triangular $p \times p$ matrix R_X , and the p -vector c_β satisfying

$$\begin{aligned}LR_{ZX} &= P\Lambda_\theta^T Z^T X \\ R_X^T R_X &= X^T X - R_{ZX}^T R_{ZX} \\ R_X^T c_\beta &= X^T y - R_{ZX}^T c_u\end{aligned}$$

so that

$$\begin{bmatrix} P^T L & 0 \\ R_{ZX}^T & R_X^T \end{bmatrix} \begin{bmatrix} L^T P & R_{ZX} \\ 0 & R_X \end{bmatrix} = \begin{bmatrix} \Lambda_\theta^T Z^T Z \Lambda_\theta + I & \Lambda_\theta^T Z^T X \\ X^T Z \Lambda_\theta & X^T X \end{bmatrix}.$$

Solving the extended PLS problem (cont'd)

- Finally we solve

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

- 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.

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 θ , is straightforward. It involves updating Λ_θ , L_θ , R_{ZX} , R_X , calculating the penalized residual sum of squares, r_θ^2 and two 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 .