

Linear mixed model implementation in lme4

Douglas Bates
Department of Statistics
University of Wisconsin – Madison

January 31, 2007

Abstract

We describe the form of the linear mixed-effects and generalized linear mixed-effects models fit by `lmer` and give details of the representation and the computational techniques used to fit such models. These techniques are illustrated on several examples.

1 A simple example

The `Rail` data set from the `nlme` package is described in ? as consisting of three measurements of the travel time of a type of sound wave on each of six sample railroad rails. We can examine the structure of these data with the `str` function

```
> str(Rail)
'data.frame':      18 obs. of  2 variables:
 $ travel: num  55 53 54 26 37 32 78 91 85 92 ...
 $ Rail  : Ord.factor w/ 6 levels "2"<"5"<"1"<"6"<...: 3 3 3 1 1 1 5 5 5 6 ...
```

Because there are only three observations on each of the rails a dotplot (Figure 1) shows the structure of the data well.

```
> print(dotplot(Rail ~ travel, Rail, xlab = "Travel time (ms)",
+             ylab = "Rail number"))
```

In building a model for these data

```
> Rail
```

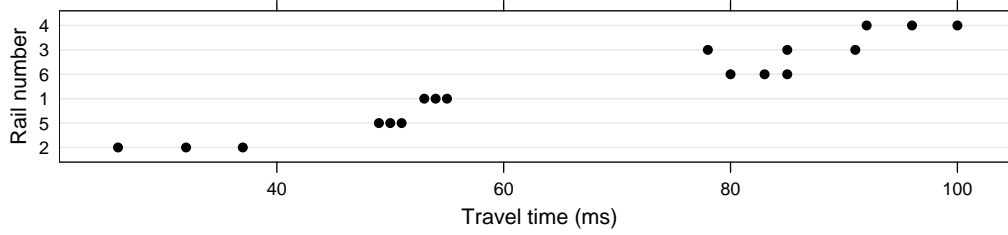


Figure 1: Travel time of sound waves in a sample of six railroad rails. There were three measurements of the travel time on each rail. The rail numbers are sorted by increasing mean travel time.

	travel	Rail
1	55	1
2	53	1
3	54	1
4	26	2
5	37	2
6	32	2
7	78	3
8	91	3
9	85	3
10	92	4
11	100	4
12	96	4
13	49	5
14	51	5
15	50	5
16	80	6
17	85	6
18	83	6

we wish to characterize a typical travel time, say μ , for the population of such railroad rails and the deviations, say $b_i, i = 1, \dots, 6$ of the individual rails from this population mean. Because these specific rails are not of interest by themselves as much as the variation in the population we model the b_i , which are called the “random effects” for the rails, as having a normal (Gaussian) distribution of the form $\mathcal{N}(0, \sigma_b^2)$. The j th measurement on the i th rail is expressed as

$$y_{ij} = \mu + b_i + \epsilon_{ij} \quad b_i \sim \mathcal{N}(0, \sigma_b^2), \epsilon_{ij} \sim \mathcal{N}(0, \sigma^2) \quad i = 1, \dots, 6 \quad j = 1, \dots, 3 \quad (1)$$

The parameters of this model are μ , σ_b^2 and σ^2 . Technically the $b_i, i = 1, \dots, 6$ are not parameters but instead are considered to be unobserved random variables for which we form “predictions” instead of “estimates”.

To express generalizations of models like (1) more conveniently we switch to a matrix/vector representation in which the 18 observations of the travel time form the response vector \mathbf{y} , the fixed-effect parameter μ forms a 1-dimensional column vector $\boldsymbol{\beta}$ and the six random effects $b_i, i = 1, \dots, 6$ form the random effects vector \mathbf{b} . The structure of the data and the values of any covariates (none are used in this model) are used to create model matrices \mathbf{X} and \mathbf{Z} .

Using these vectors and matrices and the 18-dimensional vector $\boldsymbol{\epsilon}$ that represents the per-observation noise terms the model becomes

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{b} + \boldsymbol{\epsilon}, \quad \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}), \quad \mathbf{b} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \boldsymbol{\Sigma}) \quad \text{and} \quad \mathbf{b} \perp \boldsymbol{\epsilon} \quad (2)$$

In the general form we write p for the dimension of $\boldsymbol{\beta}$, the fixed-effects parameter vector, and q for the dimension of \mathbf{b} , the vector of random effects. Thus the model matrix \mathbf{X} has dimension $n \times p$, the model matrix \mathbf{Z} has dimension $n \times q$ and the relative variance-covariance matrix, $\boldsymbol{\Sigma}$, for the random-effects has dimension $q \times q$. The symbol \perp indicates independence of random variables and \mathcal{N} denotes the multivariate normal (Gaussian) distribution.

We say that matrix $\boldsymbol{\Sigma}$ is the relative variance-covariance matrix of the random effects in the sense that it is the variance of \mathbf{b} relative to σ^2 , the scalar variance of the per-observation noise term $\boldsymbol{\epsilon}$. Although its size, q , can be very large, $\boldsymbol{\Sigma}$ is highly structured. It is symmetric, positive semi-definite and zero except for the diagonal elements and certain elements close to the diagonal.

1.1 Fitting the model and examining the results

The maximum likelihood estimates for parameters in model (1) fit to the Rail data are obtained as

```
> Rm1ML <- lmer2(travel ~ 1 + (1 | Rail), Rail, method = "ML",
+   control = list(msVerbose = 1))
0      149.289: 0.942809
1      137.531: 1.94281
2      132.389: 2.85077
3      129.942: 3.73815
4      128.945: 4.52610
5      128.629: 5.12723
6      128.566: 5.47713
7      128.560: 5.60451
8      128.560: 5.62581
9      128.560: 5.62686
10     128.560: 5.62686
```

In this fit we have set the control parameter `msVerbose` to 1 indicating that information on the progress of the iterations should be printed after every iteration. Each line gives the iteration number, the value of the deviance (negative twice the log-likelihood) and the value of the parameter s which is the standard deviation of the random effects relative to the standard deviation of the residuals.

The printed form of the model

```
> Rm1ML
Linear mixed-effects model fit by maximum likelihood
Formula: travel ~ 1 + (1 | Rail)
Data: Rail
      AIC      BIC logLik MLdeviance REMLdeviance
132.6 134.3 -64.28      128.6       122.2
Random effects:
Groups   Name      Variance Std.Dev.
Rail    (Intercept) 541.971  23.2803
Residual                17.118   4.1373
Number of obs: 18, groups: Rail, 6

Fixed effects:
              Estimate Std. Error t value
(Intercept)    66.500      9.285   7.162
```

provides additional information about the parameter estimates and some of the measures of the fit such as the log-likelihood (-64.28), the deviance for the maximum likelihood criterion (128.6), the deviance for the REML criterion (122.2), Akaike's Information Criterion (AIC= 132.6) and Schwartz's Bayesian Information Criterion (BIC= 134.3).

The model matrices \mathbf{Z} and \mathbf{X} and the negative of the response vector $-\mathbf{y}$ are stored in the `ZXyt` slot in the transposed form. Extracting the transpose of this slot

```
> t(Rm1ML@ZXyt)
18 x 8 sparse Matrix of class "dgCMatrix"

[1,] . . 1 . . . 1 -55
[2,] . . 1 . . . 1 -53
[3,] . . 1 . . . 1 -54
[4,] 1 . . . . . 1 -26
[5,] 1 . . . . . 1 -37
[6,] 1 . . . . . 1 -32
[7,] . . . . 1 . 1 -78
[8,] . . . . 1 . 1 -91
[9,] . . . . 1 . 1 -85
[10,] . . . . . 1 1 -92
[11,] . . . . . 1 1 -100
[12,] . . . . . 1 1 -96
[13,] . 1 . . . . 1 -49
[14,] . 1 . . . . 1 -51
[15,] . 1 . . . . 1 -50
```

```
[16,] . . . 1 . . 1 -80
[17,] . . . 1 . . 1 -85
[18,] . . . 1 . . 1 -83
```

The first 6 columns of this matrix are \mathbf{Z} , the seventh column is \mathbf{X} and the eighth and final column is $-\mathbf{y}$. As indicated in the display of the matrix, it is stored as a sparse matrix. The elements represented as ‘.’ are known to be zero and are not stored explicitly.

The columns of \mathbf{Z} are indicator columns (that is, the i th column has a 1 in row j if the j th observation is on rail i , otherwise it is zero) but they are not in the usual ordering. This is because the levels of the `Rail` factor have been reordered according to increasing mean response for Figure 1.

The crossproduct of the columns of this matrix are stored as a symmetric, sparse matrix in the A slot.

```
> Rm1ML@A
8 x 8 sparse Matrix of class "dsCMatrix"

[1,] 3 . . . . . 3 -95
[2,] . 3 . . . . 3 -150
[3,] . . 3 . . . 3 -162
[4,] . . . 3 . . 3 -248
[5,] . . . . 3 . 3 -254
[6,] . . . . . 3 3 -288
[7,] 3 3 3 3 3 3 18 -1197
[8,] -95 -150 -162 -248 -254 -288 -1197 89105
```

The L component of this fitted model is a Cholesky factorization of a matrix $\mathbf{A}^*(\boldsymbol{\theta})$ where $\boldsymbol{\theta}$ is a parameter vector determining $\boldsymbol{\Sigma}(\boldsymbol{\theta})$. This matrix can be factored as $\boldsymbol{\Sigma} = \mathbf{T}\mathbf{S}\mathbf{S}^\top$, where \mathbf{T} is a unit, lower triangular matrix (that is, all the elements above the diagonal are zero and all the elements on the diagonal are unity) and \mathbf{S} is a diagonal matrix with non-negative elements on the diagonal. The matrix $\mathbf{A}^*(\boldsymbol{\theta})$ is

$$\begin{aligned} \mathbf{A}^*(\boldsymbol{\theta}) &= \begin{bmatrix} \mathbf{Z}^{*\top}\mathbf{Z}^* + \mathbf{I} & \mathbf{Z}^{*\top}\mathbf{X} & -\mathbf{Z}^{*\top}\mathbf{y} \\ \mathbf{X}^\top\mathbf{Z}^* & \mathbf{X}^\top\mathbf{X} & -\mathbf{X}^\top\mathbf{y} \\ -\mathbf{y}^\top\mathbf{Z}^* & -\mathbf{y}^\top\mathbf{X} & \mathbf{y}^\top\mathbf{y} \end{bmatrix} \\ &= \begin{bmatrix} \mathbf{T}^\top\mathbf{S} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & 1 \end{bmatrix} \mathbf{A} \begin{bmatrix} \mathbf{S}\mathbf{T} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & 1 \end{bmatrix} + \begin{bmatrix} \mathbf{I} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & 0 \end{bmatrix}. \end{aligned} \quad (3)$$

For model (1) the matrices \mathbf{T} and \mathbf{S} are particularly simple, $\mathbf{T} = \mathbf{I}_6$, the 6×6 identity matrix and $\mathbf{S} = s_{1,1}\mathbf{I}_6$ where $s_{1,1} = \sigma_b/\sigma$ is the standard

deviation of the random effects relative to the standard deviation of the per-observation noise term ϵ .

The Cholesky decomposition of \mathbf{A}^* is a lower triangular sparse matrix \mathbf{L}

```
> as(Rm1ML@L, "sparseMatrix")
8 x 8 sparse Matrix of class "dtCMatrix"
```

```
[1,]  9.797  .      .      .      .      .      .      .
[2,]  .      9.797  .      .      .      .      .      .
[3,]  .      .      9.797  .      .      .      .      .
[4,]  .      .      .      9.797  .      .      .      .
[5,]  .      .      .      .      9.797  .      .      .
[6,]  .      .      .      .      .      9.797  .      .
[7,]  1.723  1.723  1.723  1.723  1.723  1.723  0.4330  .
[8,] -54.562 -86.150 -93.042 -142.435 -145.881 -165.408 -28.7977 17.06
```

As explained in later sections the matrix \mathbf{L} provides all the information needed to evaluate the ML deviance or the REML deviance as a function of $\boldsymbol{\theta}$. The components of the deviance are given in the `deviance` slot of the fitted model

```
> Rm1ML@deviance
      ML      REML      ldZ      ldX      lr2
128.560037 122.237086 27.385123 -1.673815  5.673323
```

The element labelled `ldZ` is the logarithm of the square of the determinant of the upper left 6×6 section of \mathbf{L} . This corresponds to $\log |\mathbf{Z}^{*\top} \mathbf{Z}^* + \mathbf{I}_q|$ where $\mathbf{Z}^* = \mathbf{ZTS}$. We can verify that the value 27.38292 can indeed be calculated in this way.

```
> L <- as(Rm1ML@L, "sparseMatrix")
> 2 * sum(log(diag(L)[1:6]))
[1] 27.38512
```

The `lr2` element of the `deviance` slot is the logarithm of the penalized residual sum of squares. It can be calculated as the logarithm of the square of the last diagonal element in \mathbf{L} .

```
> 2 * log(L[8, 8])
[1] 5.673323
```

For completeness we mention that the `ldX` element of the `deviance` slot is the logarithm of the product of the squares of the diagonal elements of \mathbf{L} corresponding to columns of \mathbf{X} .

```
> 2 * log(L[7, 7])
[1] -1.673815
```

This element is used in the calculation of the REML criterion.

Another slot in the fitted model object is `dims`, which contains information on the dimensions of the model and some of the characteristics of the fit.

```
> RmlML@dims
      nf      n      p      q REML glmm
      1     18      1      6      0      0
```

We can reconstruct the ML estimate of the residual variance as the penalized residual sum of squares divided by the number of observations.

```
> exp(RmlML@deviance["lr2"])/RmlML@dims["n"]
      lr2
16.16667
```

The *profiled deviance* function

$$\begin{aligned}\tilde{\mathcal{D}}(\boldsymbol{\theta}) &= \log \left| \mathbf{Z}^{*\top} \mathbf{Z}^* + \mathbf{I}_q \right| + n \log \left(1 + \frac{2\pi r^2}{n} \right) \\ &= n \left[1 + \log \left(\frac{2\pi}{n} \right) \right] + \log \left| \mathbf{Z}^{*\top} \mathbf{Z}^* + \mathbf{I}_q \right| + n \log r^2\end{aligned}\tag{4}$$

is a function of $\boldsymbol{\theta}$ only. In this case $\boldsymbol{\theta} = \sigma_1$, the relative standard deviation of the random effects, is one-dimensional. The maximum likelihood estimate (mle) of $\boldsymbol{\theta}$ minimizes the profiled deviance. The mle's of all the other parameters in the model can be derived from the estimate of this parameters.

The term $n[1 + \log(2\pi/n)]$ in (4) does not depend on $\boldsymbol{\theta}$. The other two terms, $\log \left| \mathbf{Z}^{*\top} \mathbf{Z}^* + \mathbf{I}_q \right|$ and $n \log r^2$, measure the complexity of the model and the fidelity of the fitted values to the observed data, respectively. We plot the value of each of the varying terms versus σ_1 in Figure 2.

The component $\log \left| \mathbf{S} \mathbf{Z}^\top \mathbf{Z} \mathbf{S} + \mathbf{I} \right|$ has the value 0 at $\sigma_1 = 0$ and increases as σ_1 increases. It is unbounded as $\sigma_1 \rightarrow \infty$. The component $n \log(r^2)$ has a finite value at $\sigma_1 = 1$ from which it decreases as σ_1 increases. The value at $\sigma_1 = 0$ corresponds to the residual sum of squares for the regression of \mathbf{y} on the columns of \mathbf{X} .

```
> 18 * log(deviance(lm(travel ~ 1, Rail)))
[1] 164.8714
```

As $\sigma_1 \rightarrow \infty$, $n \log(r^2)$ approaches the value corresponding to the residual sum of squares for the regression of \mathbf{y} on the columns of \mathbf{X} and \mathbf{Z} . For this model that is

```
> 18 * log(deviance(lm(travel ~ Rail, Rail)))
[1] 94.82145
```

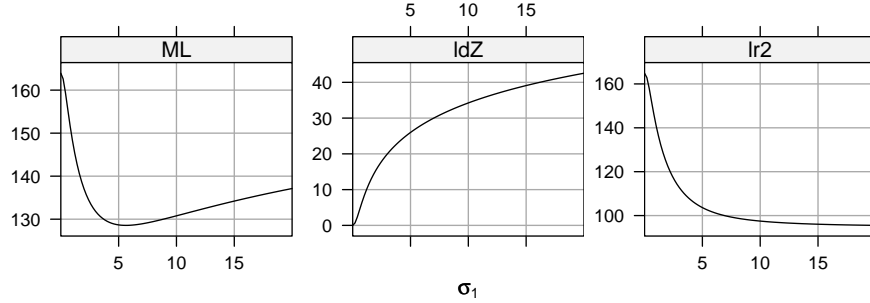


Figure 2: The profiled deviance, and those components of the profiled deviance that vary with θ , as a function of θ in model **Rm1ML** for the **Rail** data. In this model the parameter θ is the scalar σ_1 , the standard deviation of the random effects relative to the standard deviation of the per-observation noise.

2 Structure of Σ and Z

The columns of Z and the rows and columns of Σ are associated with the levels of one or more grouping factors in the data. For example, a common application of linear mixed models is the analysis of students' scores on the annual state-wide performance tests mandated by the No Child Left Behind Act. A given score is associated with a student, a teacher, a school and a school district. These could all be grouping factors in a model.

We write the grouping factors as $\mathbf{f}_i, i = 1, \dots, k$. The number of levels of the i th factor, \mathbf{f}_i , is n_i and the number of random effects associated with each level is q_i . For example, if \mathbf{f}_1 is “student” then n_1 is the number of students in the study. If we have a simple additive random effect for each student then $q_1 = 1$. If we have a random effect for both the intercept and the slope with respect to time for each student then $q_1 = 2$. The $q_i, i = 1, \dots, k$ are typically very small whereas the $n_i, i = 1, \dots, k$ can be very large.

In the statistical model we assume that random effects associated with different grouping factors are independent, which implies that Σ is block diagonal with k diagonal blocks of sizes $n_i q_i \times n_i q_i, i = 1, \dots, k$. That is

$$\Sigma = \begin{bmatrix} \Sigma_1 & \mathbf{0} & \dots & \mathbf{0} \\ \mathbf{0} & \Sigma_2 & \dots & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{0} & \mathbf{0} & \dots & \Sigma_k \end{bmatrix} \quad (5)$$

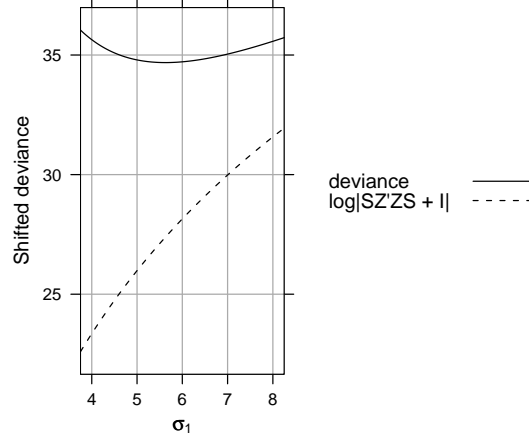


Figure 3: The part of the deviance that varies with σ_1 as a function of σ_1 near the optimum. The component $\log |\mathbf{S}\mathbf{Z}^\top \mathbf{Z}\mathbf{S} + \mathbf{I}|$ is shown at the bottom of the frame. This is the component of the deviance that increases with σ_1 . Added to this component is $n \log [r^2(\sigma_1)] - n \log [r^2(\infty)]$, the component of the deviance that decreases as σ_1 increases. Their sum is minimized at $\hat{\sigma}_1 = 5.626$.

Furthermore, random effects associated with different levels of the same grouping factor are assumed to be independent and identically distributed, which implies that Σ_i is itself block diagonal in n_i blocks and that each of these blocks is a copy of a $q_i \times q_i$ matrix $\tilde{\Sigma}_i$. That is

$$\Sigma_i = \begin{bmatrix} \tilde{\Sigma}_i & \mathbf{0} & \dots & \mathbf{0} \\ \mathbf{0} & \tilde{\Sigma}_i & \dots & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{0} & \mathbf{0} & \dots & \tilde{\Sigma}_i \end{bmatrix} = \mathbf{I}_{n_i} \otimes \tilde{\Sigma}_i \quad i = 1, \dots, k \quad (6)$$

where \otimes denotes the Kronecker product.

The condition that Σ is positive semi-definite holds if and only if the $\tilde{\Sigma}_i, i = 1, \dots, k$ are positive semi-definite. To ensure that the $\tilde{\Sigma}_i$ are positive semi-definite, we express them as

$$\tilde{\Sigma}_i = \tilde{T}_i \tilde{\mathbf{S}}_i \tilde{\mathbf{S}}_i^\top \tilde{T}_i^\top, \quad i = 1, \dots, k \quad (7)$$

where \tilde{T}_i is a $q_i \times q_i$ unit lower-triangular matrix (i.e. all the elements above

the diagonal are zero and all the diagonal elements are unity) and $\tilde{\mathbf{S}}_i$ is a $q_i \times q_i$ diagonal matrix with non-negative elements on the diagonal.

This is the “LDL” form of the Cholesky decomposition of positive semi-definite matrices except that we express the diagonal matrix \mathbf{D} , which is on the variance scale, as the square of the diagonal matrix \mathbf{S} , which is on the standard deviation scale. The profiled deviance behaves more like a quadratic on the standard deviation scale than it does on the variance scale so the use of the standard deviation scale enhances convergence.

The $n_i q_i \times n_i q_i$ matrices $\mathbf{S}_i, \mathbf{T}_i, i = 1, \dots, k$ and the $q \times q$ matrices \mathbf{S} and \mathbf{T} are defined analogously to (6) and (5). In particular,

$$\mathbf{S}_i = \mathbf{I}_{n_i} \otimes \tilde{\mathbf{S}}_i, \quad i = 1, \dots, k \quad (8)$$

$$\mathbf{T}_i = \mathbf{I}_{n_i} \otimes \tilde{\mathbf{T}}_i, \quad i = 1, \dots, k \quad (9)$$

Note that when $q_i = 1$, $\tilde{\mathbf{T}}_i = \mathbf{I}$ and hence $\mathbf{T}_i = \mathbf{I}$. Furthermore, \mathbf{S}_i is a multiple of the identity matrix in this case.

The parameter vector $\boldsymbol{\theta}_i, i = 1, \dots, k$ consists of the q_i diagonal elements of $\tilde{\mathbf{S}}_i$, which are constrained to be non-negative, followed by the $q_i(q_i - 1)/2$ elements in the strict lower triangle of $\tilde{\mathbf{T}}_i$ (in column-major ordering). These last $q_i(q_i - 1)/2$ elements are unconstrained. The $\boldsymbol{\theta}_i$ are combined as

$$\boldsymbol{\theta} = \begin{bmatrix} \boldsymbol{\theta}_1 \\ \boldsymbol{\theta}_2 \\ \vdots \\ \boldsymbol{\theta}_k \end{bmatrix}.$$

Each of the $q \times q$ matrices \mathbf{S}, \mathbf{T} and $\boldsymbol{\Sigma}$ in the decomposition $\boldsymbol{\Sigma} = \mathbf{T} \mathbf{S} \mathbf{S}^T$ is a function of $\boldsymbol{\theta}$.

As a unit triangular matrix \mathbf{T} is non-singular. That is, \mathbf{T}^{-1} exists and is easily calculated from the $\tilde{\mathbf{T}}_i^{-1}, i = 1, \dots, k$. When $\boldsymbol{\theta}$ is not on the boundary defined by the constraints, \mathbf{S} is a diagonal matrix with strictly positive elements on the diagonal, which implies that \mathbf{S}^{-1} exists and that $\boldsymbol{\Sigma}$ is non-singular with $\boldsymbol{\Sigma}^{-1} = \mathbf{T}^{-T} \mathbf{S}^{-1} \mathbf{S}^{-1} \mathbf{T}^{-1}$.

When $\boldsymbol{\theta}$ is on the boundary the matrices \mathbf{S} and $\boldsymbol{\Sigma}$ exist but are not invertible. We say that $\boldsymbol{\Sigma}$ is a *degenerate* variance-covariance matrix in the sense that one or more linear combinations of the vector \mathbf{b} are defined to have zero variance. That is, the distribution of these linear combinations is a point mass at 0.

The maximum likelihood estimates of $\boldsymbol{\theta}$ (or the restricted maximum likelihood estimates, defined below) can be located on the boundary. That is, they can correspond to a degenerate variance-covariance matrix and we must be careful to allow for this case. However, to begin we consider the non-degenerate case.

3 Methods for non-singular Σ

When $\boldsymbol{\theta}$ is not on the boundary we can define a standardized random effects vector

$$\mathbf{b}^* = \mathbf{S}^{-1}\mathbf{T}^{-1}\mathbf{b} \quad (10)$$

with the properties

$$\mathbb{E}[\mathbf{b}^*] = \mathbf{S}^{-1}\mathbf{T}^{-1}\mathbb{E}[\mathbf{b}] \quad (11)$$

$$\begin{aligned} \text{Var}[\mathbf{b}^*] &= \mathbb{E}[\mathbf{b}^*\mathbf{b}^{*\top}] \\ &= \mathbf{S}^{-1}\mathbf{T}^{-1}\text{Var}[\mathbf{b}]\mathbf{T}^{-\top}\mathbf{S}^{-1} \\ &= \sigma^2\mathbf{S}^{-1}\mathbf{T}^{-1}\Sigma\mathbf{T}^{-\top}\mathbf{S}^{-1} \\ &= \sigma^2\mathbf{S}^{-1}\mathbf{T}^{-1}\mathbf{T}\mathbf{S}\mathbf{S}^{\top}\mathbf{T}^{-\top}\mathbf{S}^{-1} \\ &= \sigma^2\mathbf{I}. \end{aligned} \quad (12)$$

Thus, the unconditional distribution of the q elements of \mathbf{b}^* is $\mathbf{b}^* \sim \mathcal{N}(\mathbf{0}, \sigma^2\mathbf{I})$, like that of the n elements of $\boldsymbol{\epsilon}$.

Obviously the transformation from \mathbf{b}^* to \mathbf{b} is

$$\mathbf{b} = \mathbf{T}\mathbf{S}\mathbf{b}^* \quad (13)$$

and the $n \times q$ model matrix for \mathbf{b}^* is

$$\mathbf{Z}^* = \mathbf{Z}\mathbf{T}\mathbf{S} \quad (14)$$

so that

$$\mathbf{Z}^*\mathbf{b}^* = \mathbf{Z}\mathbf{T}\mathbf{S}\mathbf{S}^{-1}\mathbf{T}^{-1}\mathbf{b} = \mathbf{Z}\mathbf{b}. \quad (15)$$

Notice that \mathbf{Z}^* can be evaluated even when $\boldsymbol{\theta}$ is on the boundary. Also, if we have a value of \mathbf{b}^* in such a case, we can evaluate \mathbf{b} from \mathbf{b}^* .

Given the data \mathbf{y} and values of $\boldsymbol{\theta}$ and $\boldsymbol{\beta}$, the mode of the conditional distribution of \mathbf{b}^* is the solution to a penalized least squares problem

$$\begin{aligned}\tilde{\mathbf{b}}^*(\boldsymbol{\theta}, \boldsymbol{\beta} | \mathbf{y}) &= \arg \min_{\mathbf{b}^*} \left[\|\mathbf{y} - \mathbf{X}\boldsymbol{\beta} - \mathbf{Z}^* \mathbf{b}^*\|^2 + \mathbf{b}^{*\top} \mathbf{b}^* \right] \\ &= \arg \min_{\mathbf{b}^*} \left\| \begin{bmatrix} \mathbf{y} \\ \mathbf{0} \end{bmatrix} - \begin{bmatrix} \mathbf{Z}^* & \mathbf{X} \\ \mathbf{I} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{b}^* \\ \boldsymbol{\beta} \end{bmatrix} \right\|^2.\end{aligned}\quad (16)$$

In fact, if we optimize the penalized least squares expression in (16) with respect to both \mathbf{b} and $\boldsymbol{\beta}$ we obtain the conditional estimates $\hat{\boldsymbol{\beta}}(\boldsymbol{\theta} | \mathbf{y})$ and the conditional modes $\tilde{\mathbf{b}}^*(\boldsymbol{\theta}, \hat{\boldsymbol{\beta}}(\boldsymbol{\theta} | \mathbf{y}))$ which we write as $\hat{\mathbf{b}}^*(\boldsymbol{\theta})$. That is,

$$\begin{aligned}\begin{bmatrix} \hat{\mathbf{b}}^*(\boldsymbol{\theta}) \\ \hat{\boldsymbol{\beta}}(\boldsymbol{\theta}) \end{bmatrix} &= \arg \min_{\mathbf{b}^*, \boldsymbol{\beta}} \left\| \begin{bmatrix} \mathbf{Z}^* & \mathbf{X} & -\mathbf{y} \\ \mathbf{I} & \mathbf{0} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{b}^* \\ \boldsymbol{\beta} \\ 1 \end{bmatrix} \right\|^2 \\ &= \arg \min_{\mathbf{b}^*, \boldsymbol{\beta}} \begin{bmatrix} \mathbf{b}^* \\ \boldsymbol{\beta} \\ 1 \end{bmatrix}^\top \mathbf{A}^*(\boldsymbol{\theta}) \begin{bmatrix} \mathbf{b}^* \\ \boldsymbol{\beta} \\ 1 \end{bmatrix}\end{aligned}\quad (17)$$

where the matrix $\mathbf{A}^*(\boldsymbol{\theta})$ is as shown in (3) and

$$\mathbf{A} = \begin{bmatrix} \mathbf{Z}^\top \mathbf{Z} & \mathbf{Z}^\top \mathbf{X} & -\mathbf{Z}^\top \mathbf{y} \\ \mathbf{X}^\top \mathbf{Z} & \mathbf{X}^\top \mathbf{X} & -\mathbf{X}^\top \mathbf{y} \\ -\mathbf{y}^\top \mathbf{Z} & -\mathbf{y}^\top \mathbf{X} & \mathbf{y}^\top \mathbf{y} \end{bmatrix}.\quad (18)$$

Note that \mathbf{A} does not depend upon $\boldsymbol{\theta}$. Furthermore, the nature of the model matrices \mathbf{Z} and \mathbf{X} ensures that the pattern of nonzeros in $\mathbf{A}^*(\boldsymbol{\theta})$ is the same as that in \mathbf{A} .

Let the $q \times q$ permutation matrix \mathbf{P}_Z represent a fill-reducing permutation for $\mathbf{Z}^\top \mathbf{Z}$ and \mathbf{P}_X , of size $p \times p$, represent a fill-reducing permutation for $\mathbf{X}^\top \mathbf{X}$. These could be determined, for example, using the *approximate minimal degree* (AMD) algorithm described in Davis (2006) and Davis (1996) and implemented in both the **Csparse** (Davis, 2005b) and the **CHOLMOD** (Davis, 2005a) libraries of C functions. (In many cases $\mathbf{X}^\top \mathbf{X}$ is dense, but of small dimension compared to $\mathbf{Z}^\top \mathbf{Z}$, and $\mathbf{Z}^\top \mathbf{X}$ is nearly dense so \mathbf{P}_X can be \mathbf{I}_p , the $p \times p$ identity matrix.)

Let the permutation matrix \mathbf{P} be

$$\mathbf{P} = \begin{bmatrix} \mathbf{P}_Z & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{P}_X & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & 1 \end{bmatrix}\quad (19)$$

and $\mathbf{L}(\boldsymbol{\theta})$ be the sparse Cholesky decomposition of $\mathbf{A}^*(\boldsymbol{\theta})$ relative to this permutation. That is, $\mathbf{L}(\boldsymbol{\theta})$ is a sparse lower triangular matrix with the property that

$$\mathbf{L}(\boldsymbol{\theta})\mathbf{L}(\boldsymbol{\theta})^\top = \mathbf{P}\mathbf{A}^*(\boldsymbol{\theta})\mathbf{P}^\top \quad (20)$$

For $\mathbf{L}(\boldsymbol{\theta})$ to exist we must ensure that $\mathbf{A}^*(\boldsymbol{\theta})$ is positive definite. Examination of (17) shows that this will be true if \mathbf{X} is of full column rank and \mathbf{y} does not lie in the column span of \mathbf{X} (or, in statistical terms, if we can't fit \mathbf{y} perfectly using only the fixed effects).

Let $r > 0$ be the last element on the diagonal of \mathbf{L} . Then the minimum penalized residual sum of squares in (17) is r^2 and it occurs at $\hat{\mathbf{b}}^*(\boldsymbol{\theta})$ and $\hat{\boldsymbol{\beta}}(\boldsymbol{\theta})$, the solutions to the sparse triangular system

$$\mathbf{L}(\boldsymbol{\theta})^\top \mathbf{P} \begin{bmatrix} \hat{\mathbf{b}}^*(\boldsymbol{\theta}) \\ \hat{\boldsymbol{\beta}}(\boldsymbol{\theta}) \\ 1 \end{bmatrix} = \begin{bmatrix} \mathbf{0} \\ \mathbf{0} \\ r \end{bmatrix} \quad (21)$$

(Technically we should not write the 1 in the solution; it should be an unknown. However, for \mathbf{L} lower triangular with r as the last element on the diagonal and \mathbf{P} a permutation that does not move the last row, the solution for this "unknown" will always be 1.) Furthermore, $\log|\mathbf{Z}^{*\top}\mathbf{Z} + \mathbf{I}|$ can be evaluated as the sum of the logarithms of the first q diagonal elements of $\mathbf{L}(\boldsymbol{\theta})$.

The *profiled deviance function*, $\tilde{\mathcal{D}}(\boldsymbol{\theta})$, which is negative twice the log-likelihood of model (2) evaluated at $\boldsymbol{\Sigma}(\boldsymbol{\theta})$, $\hat{\boldsymbol{\beta}}(\boldsymbol{\theta})$ and $\hat{\sigma}^2(\boldsymbol{\theta})$, can be expressed as

$$\tilde{\mathcal{D}}(\boldsymbol{\theta}) = \log|\mathbf{Z}^{*\top}\mathbf{Z} + \mathbf{I}| + n \left(1 + \log \frac{2\pi r^2}{n} \right). \quad (22)$$

Notice that it is not necessary to solve for $\hat{\boldsymbol{\beta}}(\boldsymbol{\theta})$ or $\hat{\mathbf{b}}^*(\boldsymbol{\theta})$ or $\hat{\mathbf{b}}(\boldsymbol{\theta})$ to be able to evaluate $d(\boldsymbol{\theta})$. All that is needed is to update \mathbf{A} to form \mathbf{A}^* from which the sparse Cholesky decomposition $\mathbf{L}(\boldsymbol{\theta})$ can be calculated and $\tilde{\mathcal{D}}(\boldsymbol{\theta})$ evaluated.

Once $\hat{\boldsymbol{\theta}}$ is determined we can solve for $\hat{\boldsymbol{\beta}}(\hat{\boldsymbol{\theta}})$ and $\hat{\mathbf{b}}^*(\hat{\boldsymbol{\theta}})$ using (21) and for

$$\hat{\sigma}^2(\hat{\boldsymbol{\theta}}) = \frac{r^2(\hat{\boldsymbol{\theta}})}{n}. \quad (23)$$

Furthermore, $\hat{\mathbf{b}}(\hat{\boldsymbol{\theta}}) = \mathbf{S}\mathbf{T}\hat{\mathbf{b}}^*(\hat{\boldsymbol{\theta}})$.

4 Methods for singular Σ

When θ is on the boundary, corresponding to a singular Σ , some of the columns of \mathbf{Z}^* are zero. However, the matrix \mathbf{A}^* is non-singular and elements of \mathbf{b}^* corresponding to the zeroed columns in \mathbf{Z}^* approach zero smoothly as θ approaches the boundary. Thus $r(\theta)$ and $|\mathbf{Z}^{*\top}\mathbf{Z} + \mathbf{I}|$ are well-defined, as are $\tilde{\mathcal{D}}(\theta)$ and the conditional modes $\hat{\mathbf{b}}(\theta)$.

In other words, (3) and (20) can be used to define $\tilde{\mathcal{D}}(\theta)$ whether or not θ is on the boundary.

5 REML estimates

It is common to estimate the per-observation noise variance σ^2 in a fixed-effects linear model as $\hat{\sigma}^2 = r^2/(n - p)$ where r^2 is the (unpenalized) residual sum-of-squares, n is the number of observations and p is the number of fixed-effects parameters. This is not the maximum likelihood estimate of σ^2 , which is r^2/n . It is the “restricted” or “residual” maximum likelihood (REML) estimate, which takes into account that the residual vector $\mathbf{y} - \hat{\mathbf{y}}$ is constrained to a linear subspace of dimension $n - p$ in the response space. Thus its squared length, $\|\mathbf{y} - \hat{\mathbf{y}}\|^2 = r^2$, has only $n - p$ *degrees of freedom* associated with it.

The profiled REML deviance for a linear mixed model can be expressed as

$$\tilde{\mathcal{D}}_R(\theta) = \log |\mathbf{Z}^{*\top}\mathbf{Z}^* + \mathbf{I}| + \log |\mathbf{L}_\mathbf{X}|^2 + (n - p) \left(1 + \log \frac{2\pi r^2}{n - p} \right). \quad (24)$$

References

- Tim Davis. CHOLMOD: sparse supernodal Cholesky factorization and update/downdate. <http://www.cise.ufl.edu/research/sparse/cholmod>, 2005a.
- Tim Davis. CSparse: a concise sparse matrix package. <http://www.cise.ufl.edu/research/sparse/CSparse>, 2005b.
- Tim Davis. An approximate minimal degree ordering algorithm. *SIAM J. Matrix Analysis and Applications*, 17(4):886–905, 1996.
- Timothy A. Davis. *Direct Methods for Sparse Linear Systems*. Fundamentals of Algorithms. SIAM, 2006.