

# Weighted Linear Mixed-Effects Models

Developed by Paul Bailey<sup>\*†</sup>

February 14, 2018

## Introduction

{introduction}

Obtaining population estimates for official statistics or doing social science research with hopes of a direct appeal to external validity requires that a population be sampled. That usually means multi-stage probability sampling, with groups drawn, potentially groups in groups, and then individuals. These samples can be analyzed with weighted linear models and variance estimators applied Wolter (year). Another common approach to nested models is to use hierarchical linear models, or a subset of mixed models where the units are entirely nested in groups.

While any likelihood that can be written down can be maximized, efficient calculation is always appreciated by analysts. In particular, the weighted hierarchical linear model has nested integrals that are difficult to compute efficiently. For unweighted data, Bates and Pinheiro (1998) show a method of calculation that leaves only a probability density function in the integral, which can be easily integrated as the multiplicative identity, obviating the need for numerical integration.

This paper generalizes the method in Bates and Pinheiro (1998) and the `lme4` package (Bates, Maechler, Bolker, and Walker, 2015) to the weighted case. I am unaware of a description of this and will detail the information available in existing statistical software documentation which tends to simply cite Rabe-Hesketh and not provide additional detail on how weighted estimates are formed.

This paper develops the method used to more quickly estimate linear models. These are implemented in the `WeMix` R package. After the derivation, a few example growth curve models are fit with `WeMix`.

## Model

{model}

The model of interest is a linear mixed model where units are nested inside groups, which may themselves be nested in groups, a model specified by Rabe-Hesketh and Skrondal (2006) and Rabe-Hesketh, Skrondal, and Pickles (2002) for the GLLAMM software. The advantage of fitting the model in nested levels is that survey sampling weights can be incorporated; the pseudo-likelihood,<sup>1</sup> at the top level, is an integral that sums across the top-level groups, weighting each group or unit (in a two-level model) according to sampling weights.

At the highest level, the likelihood ( $\mathcal{L}$ ) is a function of the parameters for the random-effect covariance matrix ( $\boldsymbol{\theta}$ ), the fixed-effect estimates ( $\boldsymbol{\beta}$ ), and the outcomes  $\mathbf{y}$ . This equation relates the overall likelihood  $[\mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\beta}; \mathbf{y})]$  to the integral of the likelihoods of the integrals of the top-level units  $[\mathcal{L}_j^{(L)}(\boldsymbol{\theta}, \boldsymbol{\beta}; \mathbf{y}|\mathbf{u}^{(L)})]$ , indexed by  $j$ ; the top-level unit likelihoods have a superscript ( $L$ ) to indicate they are for the  $L$ th (top) level. Here we omit reference to the fixed-effect design matrix  $\mathbf{X} \in \mathbb{R}^{n \times k_x}$  and the random-effect design matrix  $\mathbf{Z} \in \mathbb{R}^{n \times k_z}$ , but

---

<sup>\*</sup>This publication was prepared for NCES (National Center for Education Statistics) under Contract No. ED-IES-12-D-0002 with American Institutes for Research. Mention of trade names, commercial products, or organizations does not imply endorsement by the U.S. government.

<sup>†</sup>The authors would like to thank Mike Cohen for reviewing this document.

<sup>1</sup>For inverse probability of selection survey data, the likelihood is an estimated population likelihood and not an observed likelihood; the typical phrase for this is *pseudo-likelihood*, which applies to every likelihood in this document.

the likelihood is conditional on those terms as well.

$$\mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\beta}; \mathbf{y}) = \int g^{(L)}(\mathbf{u}^{(L)}; \boldsymbol{\theta}^{(L)}) \prod_j \left[ \mathcal{L}_j^{(L)}(\boldsymbol{\theta}, \boldsymbol{\beta}; \mathbf{y} | \mathbf{u}^{(L)}) \right]^{w_j^{(L)}} d\mathbf{u}^{(L)} \quad (1)$$

where the  $\mathbf{u}$  terms are the random effects, which are marginalized by integrating over them;<sup>2</sup> the  $g$  function plays a role similar to a prior, but has its variance (or covariance matrix) fit using covariance parameter vector  $\boldsymbol{\theta}$ , while  $j$  represents the indexes of all the top-level units, which have their likelihood raised to the power of their weight  $w_j^{(L)}$ . Because  $\mathbf{u}$  may be a vector—for example, if there is a random slope and intercept—the covariance matrix between the  $\mathbf{u}$  terms ( $\boldsymbol{\Sigma}^{(L)}$ ) may be diagonal (no covariance) or dense. In any case, the covariance matrix is parameterized by a vector of values  $\boldsymbol{\theta}$ ; at the  $L$ th level, the relevant elements are denoted by  $\boldsymbol{\theta}^{(L)}$ .

The conditional likelihood at each level from  $L$  to two—those above the first, or lowest level—is then recursively defined, for the  $j$ th unit, at level  $l$  ( $\mathcal{L}_j^{(l)}$ ; Rabe-Hesketh et al., 2002, eq. 3) as:

$$\mathcal{L}_j^{(l)}(\boldsymbol{\theta}, \boldsymbol{\beta}; \mathbf{y} | \mathbf{U}^{(l+1)}) = \int g^{(l)}(\mathbf{u}^{(l)}; \boldsymbol{\theta}^{(l)}) \prod_{j'} \left[ \mathcal{L}_{j'}^{(l-1)}(\boldsymbol{\theta}, \boldsymbol{\beta}; \mathbf{y} | \mathbf{U}^{(l)}) \right]^{w_{j'}^{(l-1)}} d\mathbf{u}^{(l)} \quad l = 2, \dots, L-1 \quad (2)$$

where the subscript  $j'$  that the product is indexed over indicates that the likelihood  $\mathcal{L}_{j'}^{(l-1)}(\cdot)$  is for the units of level  $l-1$  nested in unit  $j$ , and  $\mathbf{U}^{(l)}$  is all of the random effects at or above level  $l$ , so that  $\mathbf{U}^{(l)}$  includes  $\{\mathbf{u}^{(l)}, \mathbf{u}^{(l+1)}, \dots, \mathbf{u}^{(L)}\}$ . This equation reflects the nested nature of the data; a group (e.g., school), annotated as  $j$ , has an individual likelihood that is the product of the likelihoods of the units or groups (e.g., students or classrooms, annotated as  $j'$ ) nested inside of it.

In the case of a Gaussian residual, the likelihood ( $\mathcal{L}^{(1)}$ ) at the individual unit level is given by the equation

$$\mathcal{L}_i^{(1)}(\boldsymbol{\theta}, \boldsymbol{\beta}; \mathbf{y}, \mathbf{U}^{(2)}) = \frac{1}{\sigma} \phi\left(\frac{\hat{e}_i^{(1)}}{\sigma}\right) \quad (3)$$

where the subscript  $i$  is used to indicate that this is the likelihood of the  $i^{th}$  individual, the superscript (1) on  $\hat{e}_i^{(1)}$  is used to indicate that it is an unpenalized residual—where the penalized residual will be introduced in the next section— $\phi(\cdot)$  is the standard normal density function and  $\sigma$  is the residual variance (a scalar), and the residuals vector  $\hat{\mathbf{e}}^{(1)}$  represents the residuals  $\hat{e}_i^{(1)}$  for each individual and is given, in vector form, by

$$\hat{\mathbf{e}}^{(1)} = \mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}} - \sum_{l=2}^L \mathbf{Z}^{(l)} \hat{\mathbf{u}}^{(l)} \quad (4)$$

where  $\mathbf{Z}^{(l)} \in \mathbb{R}^{n \times k_z^{(l)}}$  are block matrixes that form  $\mathbf{Z}$

$$\mathbf{Z} = [\mathbf{Z}^{(2)}, \dots, \mathbf{Z}^{(L)}] \quad (5)$$

When solved with the above integrals (as in Rabe-Hesketh et al., 2002),  $\sigma$  is fit as a parameter and there is no direct equation for it.

---

<sup>2</sup>The  $\mathbf{u}$  terms are integrated out and so they are not, strictly speaking, estimated. Because of this, we do not call them empirical Bayes estimates. However, the ‘WeMix’ package does provide a value for them and, because the variance of the distribution that acts as a prior is estimated, these values could reasonably be regarded as an empirical Bayes estimate of the random effects.

## Meaning of weights

{meaning-of-w

The weights used in these models are typically not the inverse selection probabilities. This would fit the random effect variances as if the data had been imputed to a sample the size of the population.

Instead the weights are often adjusted to capture both the unequal selection probabilities and total sample size. These methods are documented in Pfeiffermann, Skinner, Holmes, Goldstein, & Rasbash (1998) and Rabe-Hesketh & Skrondal (2006), and other places.

Typically, the difference between weights that are interpretable as sampling and precision weights shows up only in variance estimation. However, I'll document later that the `lme4` and `WeMix` weights result in different estimates. This is because `WeMix` uses the formulas in this paper that weight both the units as well as the groups.

## Existing literature

{existing-lit

The SAS manual is often a source of exhaustive documentation of implementation but this author cannot find any indication that these methods are used by SAS at in the PROC GLIMMIX procedure documentation, "Pseudo-likelihood Estimation for Weighted Multilevel Models."

I would commend the authors of the Stata manuals for routinely producing extensive documentation on methods and formulas of their software. Uncharacteristically, in the case of weighted HLM they simply state, "Weighted estimation is achieved through incorporating  $w_j$  and  $w_{i|j}$  into the matrix decomposition methods detailed above to reflect replicated clusters for  $w_j$  and replicated observations within clusters for  $w_{i|j}$ . (Stata, page 534)"<sup>3</sup>

## Unweighted case

{unweighted-c

This document describes how `WeMix` uses symbolic integration that relies on a mix of both Bates and Pinheiro (1998) and the `lme4` package (Bates, Maechler, Bolker, and Walker, 2015), obviating the need for numerical integration in a weighted, nested model.

The basic model in `lme4` is of the form (Bates et al., 2015, eqs. 2 and 3):

$$(\mathbf{y}|\mathbf{U} = \mathbf{u}) \sim N(\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u}, \sigma^2\mathbf{W}^{-1}) \quad (6) \quad \text{\texttt{\{eq:lme4A\}}}$$

$$\mathbf{U} \sim N(0, \boldsymbol{\Sigma}(\boldsymbol{\theta})) \quad (7) \quad \text{\texttt{\{eq:lme4B\}}}$$

where  $N(\cdot, \cdot)$  is the multivariate normal distribution, and  $\boldsymbol{\Sigma}(\boldsymbol{\theta})$  is positive semidefinite—allowing, for example, a variance parameter to be exactly zero—that is parameterized by a vector of parameters  $\boldsymbol{\theta}$ .

The likelihood is maximized by integrating (symbolically) over  $\mathbf{U}$  (Bates et al., 2015, eqs. 20, 21, 22, and 24):

$$f_{\mathbf{y}}(\mathbf{y}; \boldsymbol{\beta}, \boldsymbol{\theta}) = \int f_{\mathbf{y}|\mathbf{U}=\mathbf{u}}(\mathbf{y}; \boldsymbol{\beta}, \boldsymbol{\theta}, \mathbf{u}) \cdot f_{\mathbf{U}}(\mathbf{u}) d\mathbf{u} \quad (8) \quad \text{\texttt{\{eq:lme4C\}}}$$

where the  $f_{\mathbf{U}}$  term is analogous to  $g(\mathbf{u}^{(l)})$  in the Rabe-Hesketh et al. (2006) formulation but is intentionally represented in a non-nested structure to allow for crossed terms.

In this document we show how `WeMix` uses a derivation similar to that in Bates et al. (2015) and Bates and Pinheiro (1998) to fit a sample-weighted mixed model, avoiding the integration necessary in GLLAMM. Comparing `WeMix` to `lmer`, the latter is more general in the sense of allowing crossed terms, while `WeMix` allows for sampling weights; unweighted, the models should agree.

The next section follows Bates et al. (2015) and shows the `lme4` and `WeMix` model without weights—the only case where they are identical. The subsequent section then shows the adaptations to the likelihood for the

<sup>3</sup>See "Methods and Formulas" for "mixed—Multilevel mixed-effects linear regression" in "[ME] Stata Multilevel mixed-effects reference manual: release 17."

sample weighted case. Notably, `lme4` has unit-level weights, but they are precision weights and not level-1 sample weights, so even when only the level-1 weights are nontrivial, the models disagree. The final section describes the application of the profile likelihood (again, following `lme4`) used to estimate linear models in `WeMix`.

Following the logic of Bates et al. (2015), the unweighted (unit weight) case simplifies eqs. 57 and 58 to

$$(\mathbf{y}|\mathbf{U} = \mathbf{u}) \sim T_1 \quad (9) \quad \{\text{eq:WeMixAnoW}\}$$

$$\mathbf{U} \sim T_2 * T_3 * \dots * T_L \quad (10) \quad \{\text{eq:WeMixBnoW}\}$$

where eqs. 59 and 60 are simplified to

$$T_1 \sim N(\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u}, \sigma^2 \mathbf{I}) \quad (11) \quad \{\text{eq:WeMixA2no}\}$$

$$T_l \sim N(0, \boldsymbol{\Sigma}_l) \quad l = 2, \dots, L \quad (12) \quad \{\text{eq:WeMixB2no}\}$$

the random-effect vector  $\mathbf{U}$  is rewritten as the product of a square root-covariance matrix  $\boldsymbol{\Lambda}$  and an *iid* normal vector  $\mathbf{v}$ :

$$\mathbf{U} = \boldsymbol{\Lambda}\mathbf{v} \quad (13) \quad \{\text{eq:Uf}\}$$

$$\mathbf{v} \sim N(0, \sigma^2 \mathbf{I}) \quad (14)$$

When  $\boldsymbol{\Lambda}$  is a *square root matrix* of  $\boldsymbol{\Sigma}$ , meaning

$$\frac{1}{\sigma^2} \boldsymbol{\Sigma} = \frac{1}{\sigma^2} \boldsymbol{\Lambda}^T \boldsymbol{\Lambda} \quad (15) \quad \{\text{eq:root}\}$$

it follows that  $\mathbf{U}$  has the distribution shown in eq. 12, but the equations can use the much easier to work with  $\mathbf{v}$ . Note that eq. 15 implies

$$\boldsymbol{\Sigma} = \begin{bmatrix} \boldsymbol{\Sigma}_{11} & 0 & \dots & 0 \\ 0 & \boldsymbol{\Sigma}_{22} & 0 & \vdots \\ \vdots & \vdots & \ddots & \vdots \\ 0 & \dots & 0 & \boldsymbol{\Sigma}_{LL} \end{bmatrix} \quad (16)$$

$$= \begin{bmatrix} \boldsymbol{\Lambda}_{11}^T & 0 & \dots & 0 \\ 0 & \boldsymbol{\Lambda}_{22}^T & 0 & \vdots \\ \vdots & \vdots & \ddots & \vdots \\ 0 & \dots & 0 & \boldsymbol{\Lambda}_{LL}^T \end{bmatrix} \begin{bmatrix} \boldsymbol{\Lambda}_{11} & 0 & \dots & 0 \\ 0 & \boldsymbol{\Lambda}_{22} & 0 & \vdots \\ \vdots & \vdots & \ddots & \vdots \\ 0 & \dots & 0 & \boldsymbol{\Lambda}_{LL} \end{bmatrix} \quad (17)$$

$$= \boldsymbol{\Lambda}^T \boldsymbol{\Lambda} \quad (18)$$

Note that  $\boldsymbol{\Lambda}$  (and thus  $\boldsymbol{\Sigma}$ ) will be parameterized by a set of parameters  $\boldsymbol{\theta}$ , so eq. 15 could be written

$$\frac{1}{\sigma^2} \boldsymbol{\Sigma}(\boldsymbol{\theta}) = \frac{1}{\sigma^2} (\boldsymbol{\Lambda}(\boldsymbol{\theta}))^T \boldsymbol{\Lambda}(\boldsymbol{\theta}) \quad (19)$$

and will be written this way to keep track of the parameters involved. These parameters vary by model. For a two-level model with a random intercept and nothing else,  $\boldsymbol{\theta}$  would be a scalar with the relative precision of the random effect. The matrix  $\boldsymbol{\Lambda}(\boldsymbol{\theta})$  would then be diagonal and of the form  $\boldsymbol{\Lambda}(\boldsymbol{\theta}) = \boldsymbol{\theta} \mathbf{I}$ , with  $\mathbf{I}$  being an identity matrix with the same number of rows and columns as there are groups (random effects). For a two-level model with a random slope and intercept, the  $\boldsymbol{\theta}$  would have length three and would parameterize the three elements of a covariance matrix. In this special case, the parameterization could be  $\boldsymbol{\Lambda}(\boldsymbol{\theta}) = \begin{bmatrix} \theta_1 \mathbf{I} & \theta_2 \mathbf{I} \\ 0 & \theta_3 \mathbf{I} \end{bmatrix}$ ; a block matrix that allows the slope and intercept term to covary with each other, within a group.

Then, the residual (penalized) sum of squares is (Bates et al., 2015, eqs. 14 and 15)

$$r^2(\boldsymbol{\theta}, \boldsymbol{\beta}, \mathbf{v}) = \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta} - \mathbf{Z}\boldsymbol{\Lambda}(\boldsymbol{\theta})\mathbf{v}\|^2 + \|\mathbf{v}\|^2 \quad (20) \quad \{\text{eq:pseudoSS}\}$$

where  $\|\mathbf{v}\|^2$  is the sum of squares for the vector  $\mathbf{v}$ ; as an equation,  $\|\mathbf{v}\|^2 = \sum v_i^2$ .

Notice that the penalty function ( $\|\mathbf{v}\|^2$ ) and the residual both are assumed to be independent identically distributed normal distributions with variance  $\sigma^2$ ; this allows for both the regression and the random effects to be estimated in one regression where the pseudo-data outcome for the random effects is a vector of zeros (Bates et al., 2015, eq. 16). This rewrites  $r$  as a sum of squares, adding a vector of zeros below  $\mathbf{y}$ —the *pseudo-data*:

$$r^2(\boldsymbol{\theta}, \boldsymbol{\beta}, \mathbf{v}) = \left\| \begin{bmatrix} \mathbf{y} \\ \mathbf{0} \end{bmatrix} - \begin{bmatrix} \mathbf{Z}\boldsymbol{\Lambda}(\boldsymbol{\theta}) & \mathbf{X} \\ \mathbf{I} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{v} \\ \boldsymbol{\beta} \end{bmatrix} \right\|^2 \quad (21) \quad \{\text{eq:WeMixR0}\}$$

Unlike Bates et al. (2015, eq. 16), we proceed by taking a QR decomposition (Trefethen and Bau, 1997) of

$$\mathbf{A} \equiv \begin{bmatrix} \mathbf{Z}\boldsymbol{\Lambda}(\boldsymbol{\theta}) & \mathbf{X} \\ \mathbf{I} & \mathbf{0} \end{bmatrix} \quad (22) \quad \{\text{eq:uwA}\}$$

Plugging eq. 22 into eq. 21 and finding the least squares solution to

$$\mathbf{A} \begin{bmatrix} \mathbf{v} \\ \boldsymbol{\beta} \end{bmatrix} = \begin{bmatrix} \mathbf{y} \\ \mathbf{0} \end{bmatrix} \quad (23) \quad \{\text{eq:lsqno1sq}\}$$

using the QR decomposition on  $\mathbf{A}$ , which rewrites  $\mathbf{A} = \mathbf{Q}\mathbf{R}$  for an orthogonal matrix  $\mathbf{Q}$  (So,  $\mathbf{Q}^T\mathbf{Q} = \mathbf{I}$ ) and an upper triangular matrix  $\mathbf{R}$ ,

$$\mathbf{Q}\mathbf{R} \begin{bmatrix} \mathbf{v} \\ \boldsymbol{\beta} \end{bmatrix} = \begin{bmatrix} \mathbf{y} \\ \mathbf{0} \end{bmatrix} \quad (24)$$

where  $\mathbf{R}$  can be written in block form as

$$\mathbf{R} = \begin{bmatrix} \mathbf{R}_{11} & \mathbf{R}_{12} \\ \mathbf{0} & \mathbf{R}_{22} \end{bmatrix} \quad (25) \quad \{\text{eq:Rblock}\}$$

in which  $\mathbf{R}_{11}$  is also upper triangular, square, and conformable with  $\mathbf{v}$ , and  $\mathbf{R}_{22}$  is similarly upper triangular, square, and conformable with  $\boldsymbol{\beta}$ , while  $\mathbf{R}_{12}$  is rectangular and (potentially) dense.

Rewriting  $r^2$  as a deviation from the least squares solution, eq. 21 becomes

$$r^2(\boldsymbol{\theta}, \boldsymbol{\beta}, \mathbf{v}) = \left\| \begin{bmatrix} \mathbf{y} \\ \mathbf{0} \end{bmatrix} - \mathbf{A} \begin{bmatrix} \mathbf{v} \\ \boldsymbol{\beta} \end{bmatrix} \right\|^2 \quad (26)$$

$$= \left\| \begin{bmatrix} \mathbf{y} \\ \mathbf{0} \end{bmatrix} - \mathbf{A} \begin{bmatrix} \mathbf{v} - \hat{\mathbf{v}} + \hat{\mathbf{v}} \\ \boldsymbol{\beta} - \hat{\boldsymbol{\beta}} + \hat{\boldsymbol{\beta}} \end{bmatrix} \right\|^2 \quad (27)$$

$$= \left\| \begin{bmatrix} \mathbf{y} \\ \mathbf{0} \end{bmatrix} - \mathbf{A} \begin{bmatrix} \hat{\mathbf{v}} \\ \hat{\boldsymbol{\beta}} \end{bmatrix} - \mathbf{A} \begin{bmatrix} \mathbf{v} - \hat{\mathbf{v}} \\ \boldsymbol{\beta} - \hat{\boldsymbol{\beta}} \end{bmatrix} \right\|^2 \quad (28)$$

Where  $\hat{\mathbf{v}}$  and  $\hat{\boldsymbol{\beta}}$  are the least squares solution to eq. 23. Using the identity that for any vector  $\mathbf{v}$  the sum of squares is also just the inner product, so  $\|\mathbf{v}\|^2 = \mathbf{v}^T\mathbf{v}$ ,

$$r^2(\boldsymbol{\theta}, \boldsymbol{\beta}, \mathbf{v}) = \left[ \begin{bmatrix} \mathbf{y} \\ \mathbf{0} \end{bmatrix} - \mathbf{A} \begin{bmatrix} \hat{\mathbf{v}} \\ \hat{\boldsymbol{\beta}} \end{bmatrix} - \mathbf{A} \begin{bmatrix} \mathbf{v} - \hat{\mathbf{v}} \\ \boldsymbol{\beta} - \hat{\boldsymbol{\beta}} \end{bmatrix} \right]^T \left[ \begin{bmatrix} \mathbf{y} \\ \mathbf{0} \end{bmatrix} - \mathbf{A} \begin{bmatrix} \hat{\mathbf{v}} \\ \hat{\boldsymbol{\beta}} \end{bmatrix} - \mathbf{A} \begin{bmatrix} \mathbf{v} - \hat{\mathbf{v}} \\ \boldsymbol{\beta} - \hat{\boldsymbol{\beta}} \end{bmatrix} \right] \quad (29) \quad \{\text{eq:finalT}\}$$

Defining  $\hat{\mathbf{e}}$  to be the penalized least squares residual,

$$\hat{\mathbf{e}} \equiv \begin{bmatrix} \mathbf{y} \\ \mathbf{0} \end{bmatrix} - \mathbf{A} \begin{bmatrix} \hat{\mathbf{v}} \\ \hat{\boldsymbol{\beta}} \end{bmatrix} \quad (30)$$

then eq. 29 can be rewritten

$$r^2(\boldsymbol{\theta}, \boldsymbol{\beta}, \mathbf{v}) = \left[ \hat{\mathbf{e}} - \mathbf{A} \begin{bmatrix} \mathbf{v} - \hat{\mathbf{v}} \\ \boldsymbol{\beta} - \hat{\boldsymbol{\beta}} \end{bmatrix} \right]^T \left[ \hat{\mathbf{e}} - \mathbf{A} \begin{bmatrix} \mathbf{v} - \hat{\mathbf{v}} \\ \boldsymbol{\beta} - \hat{\boldsymbol{\beta}} \end{bmatrix} \right] \quad (31)$$

$$= \hat{\mathbf{e}}^T \hat{\mathbf{e}} - 2\hat{\mathbf{e}}^T \left[ \mathbf{A} \begin{bmatrix} \mathbf{v} - \hat{\mathbf{v}} \\ \boldsymbol{\beta} - \hat{\boldsymbol{\beta}} \end{bmatrix} + \left[ \mathbf{A} \begin{bmatrix} \mathbf{v} - \hat{\mathbf{v}} \\ \boldsymbol{\beta} - \hat{\boldsymbol{\beta}} \end{bmatrix} \right]^T \left[ \mathbf{A} \begin{bmatrix} \mathbf{v} - \hat{\mathbf{v}} \\ \boldsymbol{\beta} - \hat{\boldsymbol{\beta}} \end{bmatrix} \right] \right] \quad (32) \quad \{\text{eq:quad}\}$$

Since  $\hat{\mathbf{e}}$ , the uniquely minimized residuals to the least squares problem is in the null of  $\mathbf{A}$ , while  $\mathbf{A}\mathbf{x}$  is in the span of  $\mathbf{A}$  for any vector  $\mathbf{x}$ , then  $\hat{\mathbf{e}}^T \mathbf{A}\mathbf{x} = 0$  for any  $\mathbf{x}$ . Thus,  $\hat{\mathbf{e}}^T \left[ \mathbf{A} \begin{bmatrix} \mathbf{v} - \hat{\mathbf{v}} \\ \boldsymbol{\beta} - \hat{\boldsymbol{\beta}} \end{bmatrix} \right] = \mathbf{0}$  and eq. 32 becomes

$$r^2(\boldsymbol{\theta}, \boldsymbol{\beta}, \mathbf{v}) = \hat{\mathbf{e}}^T \hat{\mathbf{e}} + \begin{bmatrix} \mathbf{v} - \hat{\mathbf{v}} \\ \boldsymbol{\beta} - \hat{\boldsymbol{\beta}} \end{bmatrix}^T \mathbf{A}^T \mathbf{A} \begin{bmatrix} \mathbf{v} - \hat{\mathbf{v}} \\ \boldsymbol{\beta} - \hat{\boldsymbol{\beta}} \end{bmatrix} \quad (33)$$

$$= \hat{\mathbf{e}}^T \hat{\mathbf{e}} + \begin{bmatrix} \mathbf{v} - \hat{\mathbf{v}} \\ \boldsymbol{\beta} - \hat{\boldsymbol{\beta}} \end{bmatrix}^T [\mathbf{Q}\mathbf{R}]^T [\mathbf{Q}\mathbf{R}] \begin{bmatrix} \mathbf{v} - \hat{\mathbf{v}} \\ \boldsymbol{\beta} - \hat{\boldsymbol{\beta}} \end{bmatrix} \quad (34)$$

$$= \hat{\mathbf{e}}^T \hat{\mathbf{e}} + \begin{bmatrix} \mathbf{v} - \hat{\mathbf{v}} \\ \boldsymbol{\beta} - \hat{\boldsymbol{\beta}} \end{bmatrix}^T \mathbf{R}^T \mathbf{Q}^T \mathbf{Q} \mathbf{R} \begin{bmatrix} \mathbf{v} - \hat{\mathbf{v}} \\ \boldsymbol{\beta} - \hat{\boldsymbol{\beta}} \end{bmatrix} \quad (35)$$

Then, because  $\mathbf{Q}$  is orthonormal,  $\mathbf{Q}^T = \mathbf{Q}^{-1}$  and

$$r^2(\boldsymbol{\theta}, \boldsymbol{\beta}, \mathbf{v}) = \hat{\mathbf{e}}^T \hat{\mathbf{e}} + \begin{bmatrix} \mathbf{v} - \hat{\mathbf{v}} \\ \boldsymbol{\beta} - \hat{\boldsymbol{\beta}} \end{bmatrix}^T \mathbf{R}^T \mathbf{R} \begin{bmatrix} \mathbf{v} - \hat{\mathbf{v}} \\ \boldsymbol{\beta} - \hat{\boldsymbol{\beta}} \end{bmatrix} \quad (36)$$

$$= \hat{\mathbf{e}}^T \hat{\mathbf{e}} + \left\| \mathbf{R} \begin{bmatrix} \mathbf{v} - \hat{\mathbf{v}} \\ \boldsymbol{\beta} - \hat{\boldsymbol{\beta}} \end{bmatrix} \right\|^2 \quad (37) \quad \{\text{eq:r2F}\}$$

Notice that  $\hat{\mathbf{e}}^T \hat{\mathbf{e}}$  is the value of  $r^2$  evaluated at the least squares solution (denoted by adding hats to  $\boldsymbol{\beta}$  and  $\mathbf{v}$ ), so that

$$r^2(\boldsymbol{\theta}, \hat{\boldsymbol{\beta}}, \hat{\mathbf{v}}) = \hat{\mathbf{e}}^T \hat{\mathbf{e}} \quad (38) \quad \{\text{eq:r2tau}\}$$

Plugging eqs. 38 and 25 into eq. 37 (Bates et al., 2015, eq. 19),

$$r^2(\boldsymbol{\theta}, \boldsymbol{\beta}, \mathbf{v}) = r^2(\boldsymbol{\theta}, \hat{\boldsymbol{\beta}}, \hat{\mathbf{v}}) + \left\| \begin{bmatrix} \mathbf{R}_{11} & \mathbf{R}_{12} \\ \mathbf{0} & \mathbf{R}_{22} \end{bmatrix} \begin{bmatrix} \mathbf{v} - \hat{\mathbf{v}} \\ \boldsymbol{\beta} - \hat{\boldsymbol{\beta}} \end{bmatrix} \right\|^2 \quad (39)$$

$$= r^2(\boldsymbol{\theta}, \hat{\boldsymbol{\beta}}, \hat{\mathbf{v}}) + \left\| \mathbf{R}_{11}(\mathbf{v} - \hat{\mathbf{v}}) + \mathbf{R}_{12}(\boldsymbol{\beta} - \hat{\boldsymbol{\beta}}) \right\|^2 + \left\| \mathbf{R}_{22}(\boldsymbol{\beta} - \hat{\boldsymbol{\beta}}) \right\|^2 \quad (40) \quad \{\text{eq:r2sub}\}$$

From the joint distribution of  $\mathbf{y}$  and  $\mathbf{v}$  (Bates et al., 2015, eqs. 20, 21, and 22),

$$(\mathbf{y}, \mathbf{v}) \sim N(\mathbf{y} - \mathbf{X}\boldsymbol{\beta} - \mathbf{Z}\mathbf{v}, \sigma^2 \mathbf{I}_n) * N(\mathbf{v}, \sigma^2 \mathbf{I}_{k_z}) \quad (41)$$

and the probability density function of the joint distribution of  $\mathbf{y}$  and  $\mathbf{v}$  is

$$f_{\mathbf{y}, \mathbf{v}}(\mathbf{y}, \mathbf{v}, \boldsymbol{\beta}, \boldsymbol{\theta}, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{\frac{n}{2}}} \exp \left[ \frac{-(\mathbf{y} - \mathbf{X}\boldsymbol{\beta} - \mathbf{Z}\mathbf{v})^2}{2\sigma^2} \right] \cdot \frac{1}{(2\pi\sigma^2)^{\frac{k_z}{2}}} \exp \left[ \frac{-\|\mathbf{v}\|^2}{2\sigma^2} \right] \quad (42)$$

pluggin in eq. 20

$$f_{\mathbf{y}, \mathbf{v}}(\mathbf{y}, \mathbf{v}, \boldsymbol{\beta}, \boldsymbol{\theta}, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{\frac{n}{2}}} \exp \left[ \frac{-r^2(\boldsymbol{\theta}, \boldsymbol{\beta}, \mathbf{v}) + \|\mathbf{v}\|^2}{2\sigma^2} \right] \cdot \frac{1}{(2\pi\sigma^2)^{\frac{k_z}{2}}} \exp \left[ \frac{-\|\mathbf{v}\|^2}{2\sigma^2} \right] \quad (43)$$

$$= \frac{1}{(2\pi\sigma^2)^{\frac{n+k_z}{2}}} \exp \left[ \frac{-r^2(\boldsymbol{\theta}, \boldsymbol{\beta}, \mathbf{v})}{2\sigma^2} \right] \quad (44)$$

This can then be integrated over  $\mathbf{v}$  to get the likelihood of the model (Bates et al., 2015, eqs. 25 and 26):

$$\mathcal{L}(\boldsymbol{\beta}, \boldsymbol{\theta}, \sigma^2; \mathbf{y}) = \int f_{\mathbf{y}, \mathbf{v}}(\mathbf{y}, \mathbf{v}, \boldsymbol{\beta}, \boldsymbol{\theta}) d\mathbf{v} \quad (45)$$

$$= \int \frac{1}{(2\pi\sigma^2)^{\frac{n+k_z}{2}}} \exp \frac{-r^2(\boldsymbol{\theta}, \boldsymbol{\beta}, \mathbf{v})}{2\sigma^2} d\mathbf{v} \quad (46)$$

$$= \frac{1}{(2\pi\sigma^2)^{\frac{n+k_z}{2}}} \exp \frac{-r^2(\boldsymbol{\theta}, \hat{\boldsymbol{\beta}}, \hat{\mathbf{v}}) - \left\| \mathbf{R}_{22}(\boldsymbol{\beta} - \hat{\boldsymbol{\beta}}) \right\|^2}{2\sigma^2} \int \exp \frac{-\left\| \mathbf{R}_{11}(\mathbf{v} - \hat{\mathbf{v}}) + \mathbf{R}_{12}(\boldsymbol{\beta} - \hat{\boldsymbol{\beta}}) \right\|^2}{2\sigma^2} d\mathbf{v} \quad (47) \quad \{\text{eq:lnlInt}\}$$

This can be solved with a change of variables (Bates et al., 2015, eq. 27):

$$\boldsymbol{\gamma} = \mathbf{R}_{11}(\mathbf{v} - \hat{\mathbf{v}}) + \mathbf{R}_{12}(\boldsymbol{\beta} - \hat{\boldsymbol{\beta}}) \quad (48) \quad \{\text{eq:gamma}\}$$

$$\frac{d\boldsymbol{\gamma}}{d\mathbf{v}} = \mathbf{R}_{11} \quad (49) \quad \{\text{eq:Jgamma}\}$$

Using the change of variables formula,<sup>4</sup> we add the inverse determinant of  $\mathbf{R}_{11}$  when plugging eq. 48 into 47, and using eq. 49 (Bates et al., 2015, eq. 28):

$$\mathcal{L}(\boldsymbol{\beta}, \boldsymbol{\theta}, \sigma^2; \mathbf{y}) = \frac{1}{(2\pi\sigma^2)^{\frac{n+k_z}{2}}} \exp \frac{-r^2(\boldsymbol{\theta}, \hat{\boldsymbol{\beta}}, \hat{\mathbf{v}}) - \left\| \mathbf{R}_{22}(\boldsymbol{\beta} - \hat{\boldsymbol{\beta}}) \right\|^2}{2\sigma^2} \int \exp \frac{-\|\boldsymbol{\gamma}\|^2}{2\sigma^2} |\det(\mathbf{R}_{11})|^{-1} d\boldsymbol{\gamma} \quad (50)$$

$$= \frac{1}{|\det(\mathbf{R}_{11})| (2\pi\sigma^2)^{\frac{n}{2}}} \exp \frac{-r^2(\boldsymbol{\theta}, \hat{\boldsymbol{\beta}}, \hat{\mathbf{v}}) - \left\| \mathbf{R}_{22}(\boldsymbol{\beta} - \hat{\boldsymbol{\beta}}) \right\|^2}{2\sigma^2} \left\{ \frac{1}{(2\pi\sigma^2)^{\frac{k_z}{2}}} \int \exp \frac{-\|\boldsymbol{\gamma}\|^2}{2\sigma^2} d\boldsymbol{\gamma} \right\} \quad (51)$$

and because the term in curly braces is now of a probability density function, it integrates to unity

$$\mathcal{L}(\boldsymbol{\beta}, \boldsymbol{\theta}, \sigma^2; \mathbf{y}) = \frac{1}{|\det(\mathbf{R}_{11})| (2\pi\sigma^2)^{\frac{n}{2}}} \exp \frac{-r^2(\boldsymbol{\theta}, \hat{\boldsymbol{\beta}}, \hat{\mathbf{v}}) - \left\| \mathbf{R}_{22}(\boldsymbol{\beta} - \hat{\boldsymbol{\beta}}) \right\|^2}{2\sigma^2} \quad (52) \quad \{\text{eq:WeMixL0}\}$$

Having solved the integral symbolically, this expression for the likelihood no longer has an explicit integral and has been calculated exactly, without use of numerical quadrature. This derivation is incredibly close to Bates et al. (2015), with the only modification being that we use the QR decomposition of  $\mathbf{A}$  where they used the Cholesky decomposition of  $\mathbf{A}^T \mathbf{A}$ .

The deviance function ( $D(\cdot) \equiv -2\ell(\cdot)$ ),

$$D(\boldsymbol{\beta}, \boldsymbol{\theta}, \sigma^2; \mathbf{y}) = 2\ln|\det(\mathbf{R}_{11})| + n\ln(2\pi\sigma^2) + \frac{r^2(\boldsymbol{\theta}, \hat{\boldsymbol{\beta}}, \hat{\mathbf{v}}) + \left\| \mathbf{R}_{22}(\boldsymbol{\beta} - \hat{\boldsymbol{\beta}}) \right\|^2}{\sigma^2} \quad (53) \quad \{\text{eq:WeMixD0}\}$$

Using the profile likelihood, the deviance is minimized when  $\boldsymbol{\beta} = \hat{\boldsymbol{\beta}}$  because  $\boldsymbol{\beta}$  only appears inside of a sum of squares that can be minimized (set to zero) using  $\boldsymbol{\beta} = \hat{\boldsymbol{\beta}}$  (Bates et al., 2015, eqs. 30 and 31). The profile deviance then becomes

$$D(\boldsymbol{\theta}, \sigma^2; \mathbf{y}) = 2\ln|\det(\mathbf{R}_{11})| + n\ln(2\pi\sigma^2) + \frac{r^2(\boldsymbol{\theta}, \hat{\boldsymbol{\beta}}, \hat{\mathbf{v}})}{\sigma^2} \quad (54)$$

Similarly, the value of  $\sigma^2$  can be found by taking the derivative of the profile deviance with respect to  $\sigma^2$  and setting it equal to zero. This yields

$$\widehat{\sigma^2} = \frac{r^2(\boldsymbol{\theta}, \hat{\boldsymbol{\beta}}, \hat{\mathbf{v}})}{n} \quad (55)$$

---

<sup>4</sup>See the section, “Substitution for Multiple Variables” in “Integration by Substitution” (Wikipedia, n.d.).

giving a profile deviance that is a function only of the parameter  $\theta$ :

$$D(\theta; \mathbf{y}) = 2 \ln |\det(\mathbf{R}_{11})| + n \left( \ln \left( 2\pi \widehat{\sigma^2} \right) + 1 \right) \quad (56)$$

the profiled deviance can then be numerically maximized over the parameters in  $\theta$  and the plug in estimators can then be used for  $\hat{\sigma}^2$  and  $\beta$ .

## Weighted Case

{weighted-cas}

Generally, the Rabe-Hesketh et al. (2006) model can be rewritten in a form very similar to `lme4` as

$$(\mathbf{y} | \mathbf{U} = \mathbf{u}) \sim T_1^{\omega^{(1)}} \quad (57) \quad \text{\texttt{\{eq:WeMixA\}}}$$

$$\mathbf{U} \sim T_2^{\omega^{(2)}} * T_3^{\omega^{(3)}} * \dots * T_L^{\omega^{(L)}} \quad (58) \quad \text{\texttt{\{eq:WeMixB\}}}$$

where  $T^k$  represents the convolution of  $k$  instances of  $T$  and  $*$  represents the convolution of two distributions, with a likelihood that is their product;  $w^{(l)}$  are the weights assigned to units ( $l = 1$ ) or groups ( $l > 1$ ) that are ideally<sup>5</sup> the inverse (unconditional) probability of selecting the unit or group; and the  $T$ 's have distribution

$$T_1 \sim N(\mathbf{X}\beta + \mathbf{Z}\mathbf{u}, \sigma^2 I) \quad (59) \quad \text{\texttt{\{eq:WeMixA2\}}}$$

$$T_l \sim N(0, \Sigma_{ll}) \quad l = 2, \dots, L \quad (60) \quad \text{\texttt{\{eq:WeMixB2\}}}$$

where  $\Sigma$  is block diagonal, disallowing nonzero prior correlations across levels, with the  $l$ th block being written  $\Sigma_{ll}$ .

The purpose of `WeMix` is to estimate the likelihood of the weighted model (eqs. 57 and 58). In this section we derive the weighted estimator that is analogous to the estimator used by `lme4` and is similar to Henderson (1982) but we include a deviance (and likelihood), which Henderson omits.

The first difference is in the penalized sum of squares, which now weights the residuals by the unit-level weights ( $\Omega$ ) and the random-effect penalties by the group-level weights ( $\Psi$ ):

$$r^2(\theta, \beta, \mathbf{v}) = \left\| \Omega^{\frac{1}{2}} \left( \mathbf{y} - \mathbf{X}\beta - \sum_{l=1}^L \mathbf{Z}_l \Lambda_{ll}(\theta) v_l \right) \right\|^2 + \sum_{l=2}^L \left\| (\Psi_{ll})^{\frac{1}{2}} v_l \right\|^2 \quad (61) \quad \text{\texttt{\{eq:r2sep\}}}$$

where  $\Omega$  and  $\Psi_{ll}$  are diagonal matrices with unconditional inverse probability of selection for each unit ( $\Omega$ ) or group ( $\Psi_{ll}$ ) along its diagonal. The unconditional probability that a unit or group was selected can be readily calculated as the product of a probability of its own probability of selection and the unconditional probability of the group to which it belongs.

Then, the weighted pseudo-data notation combines the two terms in eq. 61, adding a vector of *pseudo-data* to the end of the  $\mathbf{y}$  vector:

$$r^2(\theta, \beta, \mathbf{v}) = \left\| \begin{bmatrix} \Omega^{\frac{1}{2}} \mathbf{y} \\ \mathbf{0} \end{bmatrix} - \begin{bmatrix} \Omega^{\frac{1}{2}} \mathbf{Z} \Lambda(\theta) & \Omega^{\frac{1}{2}} \mathbf{X} \\ \Psi^{\frac{1}{2}} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{v} \\ \beta \end{bmatrix} \right\|^2 \quad (62) \quad \text{\texttt{\{eq:WeMixWr2\}}}$$

$$= \left\| \Omega^{\frac{1}{2}} (\mathbf{y} - \mathbf{Z} \Lambda(\theta) \mathbf{v} - \mathbf{X} \beta) \right\|^2 + \left\| \Psi^{\frac{1}{2}} \mathbf{v} \right\|^2 \quad (63)$$

where  $\mathbf{Z}$  is now a block matrix that incorporates all of the  $\mathbf{Z}$  matrices for the various levels:

$$\mathbf{Z} = [\mathbf{Z}_1 \quad \mathbf{Z}_2 \quad \dots \quad \mathbf{Z}_L] \quad (64)$$

$\Lambda(\theta)$  is a block diagonal matrix with elements  $\Lambda_{ll}(\theta)$ ,  $\Psi$  is a block diagonal matrix with elements  $\Psi_{ll}$ , and

$$\mathbf{v} = \begin{bmatrix} v_1 \\ v_2 \\ \vdots \\ v_L \end{bmatrix} \quad (65)$$

---

<sup>5</sup>This weight is only ideally this because of how weights are adjusted for total nonresponse (see Gelman 2007).



The likelihood of  $\mathbf{y}$ , conditional on  $\mathbf{v}$  is then

$$f_{\mathbf{y}|\mathbf{v}=\mathbf{v}}(\mathbf{y}, \mathbf{v}) = \prod_{i=1}^n \left[ \frac{1}{(2\pi\sigma^2)^{\frac{1}{2}}} \exp \left( -\frac{\|\mathbf{y}_i - \mathbf{X}_i\boldsymbol{\beta} - \mathbf{Z}_i\boldsymbol{\Lambda}(\boldsymbol{\theta})\mathbf{v}\|^2}{2\sigma^2} \right) \right]^{\Omega_{ii}} \quad (66)$$

$$= \prod_{i=1}^n \frac{1}{(2\pi\sigma^2)^{\frac{\Omega_{ii}}{2}}} \exp \left( -\frac{\left\| \boldsymbol{\Omega}_{ii}^{\frac{1}{2}} (\mathbf{y}_i - \mathbf{X}_i\boldsymbol{\beta} - \mathbf{Z}_i\boldsymbol{\Lambda}(\boldsymbol{\theta})\mathbf{v}) \right\|^2}{2\sigma^2} \right) \quad (67)$$

$$= \frac{1}{(2\pi\sigma^2)^{\frac{\sum_i \Omega_{ii}}{2}}} \exp \left( -\frac{\left\| \boldsymbol{\Omega}^{\frac{1}{2}} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta} - \mathbf{Z}\boldsymbol{\Lambda}(\boldsymbol{\theta})\mathbf{v}) \right\|^2}{2\sigma^2} \right) \quad (68) \quad \{\text{eq:WeMixWfyu}\}$$

And the unconditional density of  $\mathbf{v}$  is

$$f_{\mathbf{v}}(\mathbf{v}) = \prod_{j=1}^{k_z} \left[ \frac{1}{(2\pi\sigma^2)^{\frac{1}{2}}} \exp \left[ -\frac{\|\mathbf{v}_j\|^2}{2\sigma^2} \right] \right]^{\Psi_{jj}} \quad (69)$$

$$= \prod_{j=1}^{k_z} \frac{1}{(2\pi\sigma^2)^{\frac{1}{2}\|\Psi_{jj}\|_F}} \exp \left[ -\frac{\left\| \boldsymbol{\Psi}_{jj}^{\frac{1}{2}} \mathbf{v}_j \right\|^2}{2\sigma^2} \right] \quad (70)$$

$$= \frac{1}{(2\pi\sigma^2)^{\frac{1}{2}\|\boldsymbol{\Psi}_{jj}\|_F}} \exp \left[ -\frac{\left\| \boldsymbol{\Psi}^{\frac{1}{2}} \mathbf{v} \right\|^2}{2\sigma^2} \right] \quad (71) \quad \{\text{eq:WeMixWfu}\}$$

where  $\|\cdot\|_F$  is the Frobenius norm of a matrix, or the sum of all its elements.

The joint distribution of  $\mathbf{v}$  and  $\mathbf{y}$  is then the product of eqs. 68 and 71:

$$f_{\mathbf{y},\mathbf{v}}(\mathbf{y}, \mathbf{v}, \boldsymbol{\beta}, \boldsymbol{\theta}, \sigma^2) = f_{\mathbf{y}|\mathbf{v}=\mathbf{v}}(\mathbf{y}, \mathbf{v}) \cdot f_{\mathbf{v}}(\mathbf{v}) \quad (72)$$

$$= \frac{1}{(2\pi\sigma^2)^{\frac{\sum_i \Omega_{ii}}{2}}} \exp \left[ -\frac{\left\| \boldsymbol{\Omega}^{\frac{1}{2}} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta} - \mathbf{Z}\boldsymbol{\Lambda}(\boldsymbol{\theta})\mathbf{v}) \right\|^2}{2\sigma^2} \right] \cdot \frac{1}{(2\pi\sigma^2)^{\frac{\sum_j \Psi_{jj}}{2}}} \exp \left[ -\frac{\left\| \boldsymbol{\Psi}^{\frac{1}{2}} \mathbf{v} \right\|^2}{2\sigma^2} \right] \quad (73)$$

$$= \frac{1}{(2\pi\sigma^2)^{\frac{\sum_i \Omega_{ii} + \sum_j \Psi_{jj}}{2}}} \exp \left[ -\frac{\left\| \boldsymbol{\Omega}^{\frac{1}{2}} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta} - \mathbf{Z}\boldsymbol{\Lambda}(\boldsymbol{\theta})\mathbf{v}) \right\|^2 + \left\| \boldsymbol{\Psi}^{\frac{1}{2}} \mathbf{v} \right\|^2}{2\sigma^2} \right] \quad (74)$$

$$= \frac{1}{(2\pi\sigma^2)^{\frac{\sum_i \Omega_{ii} + \sum_j \Psi_{jj}}{2}}} \exp \left[ -\frac{r^2(\boldsymbol{\theta}, \boldsymbol{\beta}, \mathbf{v})}{2\sigma^2} \right] \quad (75)$$

Using the same logic for the results in eq. 40,  $r^2$  can be written as a sum of the value at the optimum ( $\hat{\boldsymbol{\beta}}$  and  $\hat{\mathbf{v}}$ ) and deviations from that:

$$f_{\mathbf{y},\mathbf{v}}(\mathbf{y}, \mathbf{v}, \boldsymbol{\beta}, \boldsymbol{\theta}, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{\frac{\sum_i \Omega_{ii} + \sum_j \Psi_{jj}}{2}}} \exp \left[ -\frac{r^2(\boldsymbol{\theta}, \hat{\boldsymbol{\beta}}, \hat{\mathbf{v}}) - \left\| \mathbf{R}_{22}(\boldsymbol{\beta} - \hat{\boldsymbol{\beta}}) \right\|^2 - \left\| \mathbf{R}_{11}(\mathbf{v} - \hat{\mathbf{v}}) + \mathbf{R}_{12}(\boldsymbol{\beta} - \hat{\boldsymbol{\beta}}) \right\|^2}{2\sigma^2} \right] \quad (76)$$

Now, finding the integral of this over  $\mathbf{v}$ ,

$$\mathcal{L}(\beta, \theta, \sigma^2; \mathbf{y}) = \int f_{\mathbf{y}, \mathbf{v}}(\mathbf{y}, \mathbf{v}, \beta, \theta, \sigma^2) d\mathbf{v} \quad (77)$$

$$= \int \frac{1}{(2\pi\sigma^2)^{\frac{\sum_i \Omega_{ii} + \sum_j \Psi_{jj}}{2}}} \exp \left[ -\frac{r^2(\theta, \hat{\beta}, \hat{\mathbf{v}}) - \|\mathbf{R}_{22}(\beta - \hat{\beta})\|^2 - \|\mathbf{R}_{11}(\mathbf{v} - \hat{\mathbf{v}}) + \mathbf{R}_{12}(\beta - \hat{\beta})\|^2}{2\sigma^2} \right] d\mathbf{v} \quad (78)$$

$$= \frac{1}{(2\pi\sigma^2)^{\frac{\sum_i \Omega_{ii} + \sum_j \Psi_{jj}}{2}}} \exp \left[ -\frac{r^2(\theta, \hat{\beta}, \hat{\mathbf{v}}) - \|\mathbf{R}_{22}(\beta - \hat{\beta})\|^2}{2\sigma^2} \right] \int \exp \left[ -\frac{\|\mathbf{R}_{11}(\mathbf{v} - \hat{\mathbf{v}}) + \mathbf{R}_{12}(\beta - \hat{\beta})\|^2}{2\sigma^2} \right] d\mathbf{v} \quad (79) \quad \{\text{eq:tmpLA}\}$$

Notice that while the unweighted integral has  $k_z$  dimensions, this weighted integral has  $\sum_j \Psi_{jj}$  dimensions—the number of (population) individuals values to integrate out.

$$\mathcal{L}(\beta, \theta, \sigma^2; \mathbf{y}) = \frac{1}{(2\pi\sigma^2)^{\frac{\sum_i \Omega_{ii}}{2}}} \exp \left[ -\frac{r^2(\theta, \hat{\beta}, \hat{\mathbf{v}}) - \|\mathbf{R}_{22}(\beta - \hat{\beta})\|^2}{2\sigma^2} \right] \quad (80)$$

$$\left\{ \frac{1}{(2\pi\sigma^2)^{\frac{\sum_j \Psi_{jj}}{2}}} \int \exp \left[ -\frac{\|\mathbf{R}_{11}(\mathbf{v} - \hat{\mathbf{v}}) + \mathbf{R}_{12}(\beta - \hat{\beta})\|^2}{2\sigma^2} \right] d\mathbf{v} \right\} \quad (81) \quad \{\text{eq:tmpLB}\}$$

However, while we know there are  $\sum_j \Psi_{jj}$  dimensions to integrate out (the number of population cases), a change of variables must maintain the dimensionality of the integration, so it is not clear how to proceed. Instead, we name the term inside the integral  $\alpha$  and use a different methodology to derive its value. Then,

$$\mathcal{L}(\beta, \theta, \sigma^2; \mathbf{y}) = \alpha(\theta; \Omega, \Psi) \frac{1}{(2\pi\sigma^2)^{\frac{\sum_i \Omega_{ii}}{2}}} \exp \left[ -\frac{r^2(\theta, \hat{\beta}, \hat{\mathbf{v}}) - \|\mathbf{R}_{22}(\beta - \hat{\beta})\|^2}{2\sigma^2} \right] \quad (82) \quad \{\text{eq:tmpL2}\}$$

where  $\alpha$  is a constant for a fixed  $\theta$  and set of weights  $\Omega$  and  $\Psi$ .

While these formulas allow for estimation of a likelihood function that allows for estimation of  $\hat{\beta}$  and  $\hat{\sigma}^2$  via profiling, they do not depend on  $\alpha$  because  $\theta$  appears as a parameter of  $\alpha$ ; optimizing the log-likelihood with respect to  $\theta$  requires all of the terms in the log-likelihood to be calculated, including  $\alpha$ .

## Calculation of $\alpha$

Bates and Pinheiro (1998) offer an unweighted method of calculating  $\alpha$  that we extend here to admit weighting. The essential insight of Bates and Pinheiro is that the variables must be remapped, per group, using an orthogonal transform that separates out the  $q_l$  random effects associated with group  $g$  at level  $l$ . In what follows we describe a three-level case, but the methods readily generalize to the  $L$ -level case.

The integral is slightly re-expressed using  $\mathbf{u}$  instead of  $\mathbf{v}$ , but instead of using  $\mathbf{\Lambda}$  it uses the  $\mathbf{\Delta}$  matrix, which is an individual block of  $\mathbf{\Lambda}$ , so  $\mathbf{\Delta}$  is defined, implicitly, by

$$\mathbf{\Lambda}(\theta) = \mathbf{\Delta}(\theta) \otimes \mathbf{I} \quad (83)$$

where  $\otimes$  is the Kronecker product.

Using this notation, the likelihood is then given by

$$\mathcal{L} = \int \prod_g \left[ \int f_{y|U}(\mathbf{y}, \boldsymbol{\theta}, \mathbf{u}_{2g}, \dots, \mathbf{v}_{Lg}) f_U(\boldsymbol{\theta}, \mathbf{u}_{2g}) d\mathbf{u}_{2g} \right] f_U(\boldsymbol{\theta}, \mathbf{u}_{3g}) d\mathbf{u}_{3g} \quad (84) \quad \{\text{eq:intalpha}\}$$

$$\propto \int \prod_g \left[ \int \exp \left[ \left\| \Omega_{gg}^{\frac{1}{2}}(\mathbf{y}_g - \mathbf{X}_g \boldsymbol{\beta} - \mathbf{Z}_g \mathbf{u}) \right\|^2 + \left\| \mathbf{u}_{2g}^T \boldsymbol{\Delta}^T \boldsymbol{\Delta} \mathbf{u}_{2g} \right\|^2 \right] d\mathbf{u}_{2g} \right] f_U(\boldsymbol{\theta}, \mathbf{u}_{3g}) d\mathbf{u}_{3g} \quad (85) \quad \{\text{eq:intalpha2}\}$$

and iteratively rewritten to symbolically integrate out the lowest level random effects, starting with level 2 and increasing until there are no remaining integrals. When this is done, we will note the change of variable associated with a weighted model and use that for  $\alpha$  in eq. 82. Because of that, the portions unrelated to the change of variable were dropped from relation 85. Thus, notice that this goal is just for units in a particular group ( $g$ ) at level 2. These results of several integrals have to be combined to solve the integral across all groups and calculate  $\alpha$ . The value of  $\alpha$  will be calculated by level and then summed; for a three-level model:

$$\alpha = \alpha_2 + \alpha_3 \quad (86)$$

While the residual sum of squares ( $r^2$ ) has been, up until now, treated for the entire data set, the notion is to think of the residual sum of squares for just group  $g$  (at level  $l$ ). Bates and Pinheiro (1998) also use the notion of  $r_g^2$ , or the  $r^2$  contribution for group  $g$ , which is defined as

$$r_g^2(\boldsymbol{\theta}, \boldsymbol{\beta}, \mathbf{v}) = \left\| \begin{bmatrix} \Omega_{gg}^{\frac{1}{2}} \mathbf{y}_g \\ \mathbf{0} \end{bmatrix} - \begin{bmatrix} \Omega_{gg}^{\frac{1}{2}} \mathbf{Z}_g & \Omega_{gg}^{\frac{1}{2}} \mathbf{X}_g \\ \boldsymbol{\Delta}_l' & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{u}_g \\ \boldsymbol{\beta} \end{bmatrix} \right\|^2 \quad (87)$$

where  $\mathbf{y}_g$ ,  $\mathbf{v}_g$ ,  $\mathbf{X}_g$ ,  $\mathbf{Z}_g$ , and  $\Omega_{gg}$  are the rows of the  $\mathbf{y}$  vector, the  $\mathbf{v}$  vector,  $\mathbf{X}$  matrix,  $\mathbf{Z}$  matrix, and  $\Omega$  matrix that are associated with group  $g$ , respectively; while  $\boldsymbol{\Delta}(\boldsymbol{\theta})$  is the full  $\boldsymbol{\Delta}(\boldsymbol{\theta})$  matrix,  $\boldsymbol{\Delta}_l'$  is a block matrix:

$$\boldsymbol{\Delta}_l' \equiv [\boldsymbol{\Delta}_l \quad \mathbf{0}] \quad (88)$$

where  $\boldsymbol{\Delta}_l$  is the portion of  $\boldsymbol{\Delta}$  associated with level  $l$ , and the  $\mathbf{0}$  matrix is entirely zeros and conforms to the portion of  $\mathbf{u}_g$  that is associated with level 3 of the model.

Expanding  $\mathbf{Z}_g$  gives

$$r_g^2(\boldsymbol{\theta}, \boldsymbol{\beta}, \mathbf{v}) = \left\| \begin{bmatrix} \Omega_{gg}^{\frac{1}{2}} \mathbf{y}_g \\ \mathbf{0} \end{bmatrix} - \begin{bmatrix} \Omega_{gg}^{\frac{1}{2}} \mathbf{Z}_{2g} & \Omega_{gg}^{\frac{1}{2}} \mathbf{Z}_{3g} & \Omega_{gg}^{\frac{1}{2}} \mathbf{X}_g \\ \boldsymbol{\Delta}_2 & \mathbf{0} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{v}_{2g} \\ \mathbf{v}_{3g} \\ \boldsymbol{\beta} \end{bmatrix} \right\|^2 \quad (89) \quad \{\text{eq:hugerg}\}$$

Starting at level 2, a change of variables is chosen via the (full) QR decomposition,

$$\begin{bmatrix} \Omega_{gg}^{\frac{1}{2}} \mathbf{Z}_{2g} \\ \boldsymbol{\Delta}_2 \end{bmatrix} = \mathbf{Q}_{2g} \begin{bmatrix} \mathbf{R}_{12g} \\ \mathbf{0} \end{bmatrix} \quad (90) \quad \{\text{eq:QR0}\}$$

where the subscript on  $\mathbf{R}_{12g}$  indicates that it is the top submatrix (1), at level 2, and for group  $g$ . Notice that the blocks are different shapes on the left- and right-hand sides of eq. 90; on the left-hand side, the top block ( $\Omega_{gg}^{\frac{1}{2}} \mathbf{Z}_{2g}$ ) has as many rows as there are observations in group  $g$  and the bottom block ( $\boldsymbol{\Delta}_2$ ) has as many rows as there are random effects at level 2, while the right-hand side is flipped. The top block ( $\mathbf{R}_{12g}$ ) has as many rows as there are random effects at level 2, while the bottom block ( $\mathbf{0}$ ) has as many rows as there are observations in group  $g$ . The reason for this is that the change makes the top block on the right-hand side a square, upper triangular matrix, which will allow the change of variables to proceed.

Because  $\mathbf{Q}_{2g}$  is orthogonal by construction,  $\mathbf{Q}_{2g}^T \mathbf{Q}_{2g} = \mathbf{I}$ , one can freely premultiply the inside of the sum of squares in eq. 89 by  $\mathbf{Q}_{2g}^T$  without changing the sum of squares:<sup>6</sup>

$$r_g^2(\boldsymbol{\theta}, \boldsymbol{\beta}, \mathbf{v}) = \left\| \mathbf{Q}_{2g}^T \left\{ \begin{bmatrix} \Omega_{gg}^{\frac{1}{2}} \mathbf{y}_g \\ \mathbf{0} \end{bmatrix} - \begin{bmatrix} \Omega_{gg}^{\frac{1}{2}} \mathbf{Z}_{2g} & \Omega_{gg}^{\frac{1}{2}} \mathbf{Z}_{3g} & \Omega_{gg}^{\frac{1}{2}} \mathbf{X}_g \\ \boldsymbol{\Delta}_2 & \mathbf{0} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{v}_{2g} \\ \mathbf{v}_{3g} \\ \boldsymbol{\beta} \end{bmatrix} \right\} \right\|^2 \quad (91) \quad \{\text{eq:Qz}\}$$

<sup>6</sup>Recalling  $\|\mathbf{x}\|^2 = \mathbf{x}^T \mathbf{x}$ , so that, for orthogonal matrix  $\mathbf{Q}$ , it is easy to see  $\|\mathbf{Q}\mathbf{x}\|^2 = \mathbf{x}^T \mathbf{Q}^T \mathbf{Q} \mathbf{x} = \mathbf{x}^T \mathbf{x}$ .

Then, defining

$$\mathbf{Q}_{2g}^T \begin{bmatrix} \Omega_{gg}^{\frac{1}{2}} \mathbf{y}_g \\ \mathbf{0} \end{bmatrix} \equiv \begin{bmatrix} \mathbf{R}_{1yg} \\ \mathbf{R}_{2yg} \end{bmatrix} \quad \mathbf{Q}_{2g}^T \begin{bmatrix} \Omega_{gg}^{\frac{1}{2}} \mathbf{Z}_{3g} \\ \mathbf{0} \end{bmatrix} \equiv \begin{bmatrix} \mathbf{R}_{13g} \\ \mathbf{R}_{23g} \end{bmatrix} \quad \mathbf{Q}_{2g}^T \begin{bmatrix} \Omega_{gg}^{\frac{1}{2}} \mathbf{X}_g \\ \mathbf{0} \end{bmatrix} \equiv \begin{bmatrix} \mathbf{R}_{1Xg} \\ \mathbf{R}_{2Xg} \end{bmatrix} \quad (92) \quad \{\text{eq:Qtdef}\}$$

similar to eq. 90, and with the same dimensions, the blocks are of different shapes on the left and right of the equation.

Multiplying the  $\mathbf{Q}_{2g}^T$  through and plugging the eq. 92 equations into eq. 91,

$$r_g^2(\boldsymbol{\theta}, \boldsymbol{\beta}, \mathbf{u}) = \left\| \begin{bmatrix} \mathbf{R}_{1yg} \\ \mathbf{R}_{2yg} \end{bmatrix} - \begin{bmatrix} \mathbf{R}_{12g} & \mathbf{R}_{13g} & \mathbf{R}_{1Xg} \\ \mathbf{0} & \mathbf{R}_{23g} & \mathbf{R}_{2Xg} \end{bmatrix} \begin{bmatrix} \mathbf{u}_{2g} \\ \mathbf{u}_{3g} \\ \boldsymbol{\beta} \end{bmatrix} \right\|^2 \quad (93)$$

it is now possible to simplify the integral, do a change of variables and integrate out level 2 random effects for group  $g$ , and solve the first integral in eq. 85 symbolically. In particular, this rewriting of the terms means there are  $q_2$  terms with  $\mathbf{u}_{2g}$  in them and  $n_g - q_2$  terms that are redefined to be orthogonal to those two terms. It is this orthogonality that allows the other terms to be removed from the integral. So, since we have just shown

$$\|\Omega_{gg}^{\frac{1}{2}}(\mathbf{y}_g - \mathbf{X}_g \boldsymbol{\beta} - \mathbf{Z}_g \boldsymbol{\Lambda}_u(\boldsymbol{\theta}) \mathbf{u})\|^2 + \|\mathbf{u}\|^2 = \left\| \begin{bmatrix} \mathbf{R}_{1yg} \\ \mathbf{R}_{2yg} \end{bmatrix} - \begin{bmatrix} \mathbf{R}_{12g} & \mathbf{R}_{13g} & \mathbf{R}_{1Xg} \\ \mathbf{0} & \mathbf{R}_{23g} & \mathbf{R}_{2Xg} \end{bmatrix} \begin{bmatrix} \mathbf{u}_{2g} \\ \mathbf{u}_{3g} \\ \boldsymbol{\beta} \end{bmatrix} \right\|^2 \quad (94)$$

$$= \|\mathbf{R}_{1yg} - \mathbf{R}_{12g} \mathbf{u}_{2g} - \mathbf{R}_{13g} \mathbf{u}_{3g} - \mathbf{R}_{1Xg} \boldsymbol{\beta}_g\|^2 + \|\mathbf{R}_{2yg} - \mathbf{R}_{23g} \mathbf{u}_{3g} - \mathbf{R}_{2Xg} \boldsymbol{\beta}_g\|^2 \quad (95)$$

which can now be substituted into eq. 85, allowing a change of variables:

$$\gamma_{2g} = \mathbf{R}_{1yg} - \mathbf{R}_{12g} \mathbf{u}_{2g} - \mathbf{R}_{13g} \mathbf{u}_{3g} - \mathbf{R}_{1Xg} \boldsymbol{\beta}_g \quad (96)$$

$$\frac{d\gamma_{2g}}{d\mathbf{u}_{2g}} = \mathbf{R}_{12g} \quad (97)$$

The value of  $\alpha_2$  in the unweighted case is now clear:

$$\alpha_{2u} = \sum_{g=1}^{n_2} |\det(\mathbf{R}_{12g})|^{-1} \quad (98)$$

where  $\alpha_{2u}$  is the unweighted alpha. This formula can be weighted simply by applying the replicate weights to the individual terms:

$$\alpha_2 = \sum_{g=1}^{n_2} \Psi_{gg} |\det(\mathbf{R}_{12g})|^{-1} \quad (99)$$

However, for the three-level model, the likelihood still has level 3 integrals. The level 3 integral can also be removed. We cannot restart this process with the original  $\mathbf{Z}$  and  $\mathbf{X}$  matrices and other components because they change with the components inside the level 2 integral. However, the  $\mathbf{R}_{2**}$  matrices are the portions of the higher level  $\mathbf{u}_{3g}$  that were not integrated out, and they can be used independent of  $\mathbf{u}_{2g}$ .

We continue on with these remapped variables, starting the unweighted case (only at level 2), and now using  $g'$  to indicate that this is a different group (at level 3), with  $n_{g'}$  subgroups in it. Each group (labeled  $i$ ) contributes an outcomes matrix  $\mathbf{R}_{2yi}$ , a matrix per level 2 group  $\mathbf{R}_{23i}$  and a matrix per fixed effects regressor  $\mathbf{X}_{2Xi}$ , for  $i = 1, \dots, n_{g'}$ . Combining these, the residual sum of squares at level 3 for the group  $g'$  is

$$r_{g'}^2(\boldsymbol{\theta}, \boldsymbol{\beta}, \mathbf{v}) = \left\| \begin{bmatrix} \mathbf{R}_{2y1} \\ \vdots \\ \mathbf{R}_{2yn_{g'}} \\ \mathbf{0} \end{bmatrix} - \begin{bmatrix} \mathbf{R}_{231} & \mathbf{R}_{2X1} \\ \vdots & \vdots \\ \mathbf{R}_{23n_{g'}} & \mathbf{R}_{2Xn_{g'}} \\ \boldsymbol{\Delta}_3 & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{u}_{3g} \\ \boldsymbol{\beta} \end{bmatrix} \right\|^2 \quad (100) \quad \{\text{eq:RunW}\}$$

Following the example at level 2 above, the QR decomposition is then used:

$$\begin{bmatrix} \mathbf{R}_{231} \\ \vdots \\ \mathbf{R}_{23n_{g'}} \\ \mathbf{\Delta}_3 \end{bmatrix} = \mathbf{Q}_{3g',u} \begin{bmatrix} \mathbf{R}_{13g',u} \\ \mathbf{0} \end{bmatrix} \quad (101)$$

where a subscript  $u$  is used to indicate that  $\mathbf{Q}_{3g',u}$  and  $\mathbf{R}_{13g',u}$  are unweighted. The remaining steps are then identical to the level 2 case;  $\mathbf{R}_{13g',u}$  is used as the change of variables, so that  $\alpha_{3u} = \sum_{g'=1}^{n_2} |\det(\mathbf{R}_{13g',u})|^{-1}$ , and the other  $\mathbf{R}_{3**}$  matrices can be used to integrate out level-4 cases and so on.

When there are level 2 conditional weights—non-unit probabilities of selection for the groups at level 2, conditional on the selection of the level 3 unit—each matrix in eq. 100 could be replicated  $\Psi_{ii}$  times. Equivalently, each matrix can be weighted by the conditional probability of selection, so that eq. 100 becomes

$$r_{g'}^2(\boldsymbol{\theta}, \boldsymbol{\beta}, \mathbf{v}) = \left\| \begin{bmatrix} \Psi_{11}^{\frac{1}{2}} \mathbf{R}_{2y1} \\ \vdots \\ \Psi_{g'g'}^{\frac{1}{2}} \mathbf{R}_{2yn_{g'}} \\ \mathbf{0} \end{bmatrix} - \begin{bmatrix} \Psi_{11}^{\frac{1}{2}} \mathbf{R}_{231} & \Psi_{11}^{\frac{1}{2}} \mathbf{R}_{2X1} \\ \vdots & \vdots \\ \Psi_{g'g'}^{\frac{1}{2}} \mathbf{R}_{23n_{g'}} & \Psi_{g'g'}^{\frac{1}{2}} \mathbf{R}_{2Xn_{g'}} \\ \mathbf{\Delta}_3 & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{u}_{3g} \\ \boldsymbol{\beta} \end{bmatrix} \right\|^2 \quad (102) \quad \{\text{eq:R}\}$$

This change leads to the same QR decomposition as the replicated case:

$$\begin{bmatrix} \Psi_{11}^{\frac{1}{2}} \mathbf{R}_{231} \\ \vdots \\ \Psi_{g'g'}^{\frac{1}{2}} \mathbf{R}_{23n_{g'}} \\ \mathbf{\Delta}_3 \end{bmatrix} = \mathbf{Q}_{3g'} \begin{bmatrix} \mathbf{R}_{13g'} \\ \mathbf{0} \end{bmatrix} \quad (103) \quad \{\text{eq:QRg3}\}$$

The weighted value of  $\alpha_3$  for the third level is then

$$\alpha_3 = \sum_{g'=1}^{n_3} \Psi_{g'g'} |\mathbf{R}_{13g'}|^{-1} \quad (104)$$

## Implementation Details

The actual implementation of the calculation is slightly different from what is above. First, when the QR decomposition is taken (eqs. 90 and 103), it is possible to stop the decomposition as is described in Bates and Pinheiro (1998). It is also possible to continue the QR decomposition for the other levels of  $\mathbf{Z}$ . Using the  $\mathbf{QR}$  on the entire matrix still results in an orthogonal component in the submatrices, and so meets the goals of the decomposition while obviating the need to form the  $\mathbf{Q}$  matrix explicitly.

Also note that, in the above, the value of  $r^2$  was never used, so the components relating to  $\mathbf{y}$  and  $\mathbf{X}$  need not be formed.

## Estimation

**{estimation}**

Continuing to follow **lme4**, the estimation uses the profile likelihood. Since  $\boldsymbol{\beta}$  appears only in the final term in quadratic form, it is immediately evident that the maximum likelihood estimator (MLE) for  $\boldsymbol{\beta}$  is  $\hat{\boldsymbol{\beta}}$ , making eq. 82 profile to

$$D(\boldsymbol{\theta}, \sigma^2; \mathbf{y}) = 2 \ln(\alpha(\boldsymbol{\theta}; \Psi, \Omega)) + \left( \sum_i \Omega_{ii} \right) \ln(2\pi\sigma^2) + \frac{\mathbf{r}^2(\boldsymbol{\theta}, \hat{\boldsymbol{\beta}}, \hat{\mathbf{v}})}{\sigma^2} \quad (105) \quad \{\text{eq:WeMixPB}\}$$

Then, the value of  $\sigma^2$  can also be profiled out by taking the derivative of the deviance with respect to  $\sigma^2$  and setting it equal to zero (Bates et al., 2015, eq. 32):

$$0 = \frac{\sum_i \Omega_{ii}}{\widehat{\sigma^2}} - \frac{r^2(\boldsymbol{\theta}, \hat{\boldsymbol{\beta}}, \hat{\mathbf{v}})}{\widehat{\sigma^4}} \quad (106)$$

rearranging

$$\widehat{\sigma^2} = \frac{r^2(\boldsymbol{\theta}, \hat{\boldsymbol{\beta}}, \hat{\mathbf{v}})}{\sum_i \Omega_{ii}} \quad (107) \quad \text{\texttt{eq:WeMixMaxs}}$$

Eq. 107 can then be plugged into eq. 105 to give

$$D(\boldsymbol{\theta}; \mathbf{y}) = 2 \ln(\alpha(\boldsymbol{\theta}; \Psi, \Omega)) + \left( \sum_i \Omega_{ii} \right) \left[ \ln(2\pi\widehat{\sigma^2}) + 1 \right] \quad (108) \quad \text{\texttt{eq:WeMixP}}$$

This function is then minimized numerically with respect to  $\boldsymbol{\theta}$ , using the profile estimates for  $\boldsymbol{\beta}$  and  $\mathbf{v}$  (eq. 62) and  $\widehat{\sigma^2}$  (eq. 108).

The estimated values are then the  $\boldsymbol{\theta}$  that maximize eq. 108, the  $\sigma^2$  value from eq. 107, and the  $\boldsymbol{\beta}$  values from solving the system of equations in eq. 62.

## Variance Estimation \texttt{variance-est}

The inverse Hessian of  $\boldsymbol{\beta}$  is given by (Bates et al., 2015, eq. 54):

$$\widehat{\text{Var}}(\boldsymbol{\beta}) = \widehat{\sigma^2} \mathbf{R}_{22}^{-1} (\mathbf{R}_{22}^T)^{-1} \quad (109)$$

with  $\mathbf{R}_{22}$  coming from eq. 82. This variance estimator assumes that the weights are information weights and so is inappropriate for survey weights.

A robust (sandwich) variance estimator is given by (Binder, 1983) is appropriate:

$$\left( \widehat{\sigma^2} \mathbf{R}_{22}^T \right)^{-1} \mathbf{J} \left( \widehat{\sigma^2} \mathbf{R}_{22}^T \right)^{-1} \quad (110) \quad \text{\texttt{eq:sandwich}}$$

where  $\mathbf{J}$  is the sum of outer products of the Jacobian matrix

$$\mathbf{J} = \frac{n_L}{n_L - 1} \sum_{g=1}^{n_L} \frac{\partial(\ell_g)}{\partial \boldsymbol{\beta}} \quad (111)$$

where  $n_L$  is the number of level- $L$  (top-level) groups,  $g$  indexes level- $L$  groups, and  $\ell_g$  is the log-likelihood for group  $g$  and all groups and units nested inside of  $g$ . The log-likelihood of the full model is

$$\ell(\boldsymbol{\beta}, \boldsymbol{\theta}, \sigma^2; \mathbf{y}) = \ln[\alpha(\boldsymbol{\theta}; \Omega, \Psi)] - \frac{\sum_i \Omega_{ii}}{2} \ln(2\pi\sigma^2) - \frac{r^2(\boldsymbol{\theta}, \hat{\boldsymbol{\beta}}, \hat{\mathbf{v}})}{2\sigma^2} - \frac{\left\| \mathbf{R}_{22}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) \right\|^2}{2\sigma^2} \quad (112) \quad \text{\texttt{eq:Vln1}}$$

where we are allowing  $\boldsymbol{\beta}$  and  $\boldsymbol{\theta}$  to vary but are fixing  $\sigma^2$  at the estimated value of  $\hat{\sigma^2}$ . This could have been annotated by making  $\hat{\boldsymbol{\beta}}(\boldsymbol{\theta})$  because  $\hat{\boldsymbol{\beta}}$  is the estimated value conditional on  $\boldsymbol{\theta}$  and appears in the equation separate from the value of  $\boldsymbol{\beta}$ , but that is not shown here.

While it would be convenient if eq. 112 could be directly broken up into a portion attributable to each group, and some encouragement appears when the first three terms can be, the final term has dependencies across multiple groups. A distinct likelihood is needed that depends only on the data in that group. This is achieved by noting that data for a particular group is also valid data for a mixed model of the same type as the global

mixed model, and so eq. 112 can be used on a single group's data to get the group log-likelihood; thus a group log-likelihood can be written using the notion of the fitted value of  $\beta$  in the group ( $\hat{\beta}_g$ )

$$\ell_g(\beta, \theta, \sigma^2; \mathbf{y}_g) = \ln[\alpha_g(\theta; \Omega, \Psi)] - \frac{\sum_{i \in g} \Omega_{ii}}{2} \ln(2\pi\sigma^2) - \frac{r^2(\theta, \hat{\beta}_g, \hat{v}_g)}{2\sigma^2} - \frac{\|R_{22g}(\hat{\beta}_g - \beta)\|^2}{2\sigma^2} \quad (113) \quad \{\text{eq:Vlnlg}\}$$

where  $\alpha_g$  is the  $\alpha$  term for group  $g$  and any groups nested in it, the sum for  $\Omega$  is just over  $i$  terms associated with group  $g$ ,  $\hat{\beta}$  and  $\hat{v}$  are the values fitted only on this group, and  $R_{22g}$  is the result of a QR on a version of  $A$  performed on just data ( $X$ ,  $Z$ ,  $y$ ,  $\Psi$ , and  $\Omega$ ) associated with group  $g$ , while the values of  $\sigma^2$ ,  $\beta$ , and  $\theta$  are the values from the value the function is being evaluated at globally. Then,

$$\ell(\beta, \theta, \sigma^2; \mathbf{y}) = \sum_g \ell_g(\beta, \theta, \sigma^2; \mathbf{y}_g) \quad (114)$$

A few notes are required at this point on how, exactly, this is calculated in degenerate cases. When the matrix  $A_g$  is singular for a group (e.g., when there is only one unit in the group), then the inestimable values of  $\beta_g$  are set to zero when forming  $\hat{\beta}_g - \beta$ . Similarly,  $R_{22g}$  may not have enough rows to form an upper-triangular matrix. The portion that can be formed (including all columns) is then used—this does not affect the computation of  $R_{22g}(\hat{\beta}_g - \beta)$ .

Using these formulas the Jacobian matrix can now be calculated numerically and the robust variance estimator formed with eq. 110.

So far this section has regarded only  $\beta$  but similar methods apply to the estimation of the variance of the random effect variance estimates ( $\theta$  and  $\sigma$ ). These variance terms have their variance estimated assuming that they are uncorrelated with the  $\beta$  terms. At each level the variance is calculated, including a term for  $\sigma$ , as

$$\text{Var}(\theta, \sigma) = (-H)^{-1} J_{\theta, \sigma} (-H)^{-1} \quad (115) \quad \{\text{varthetaeta}\}$$

where  $H$  is the Hessian of the likelihood (eq. 112) with respect to  $\theta$  and  $\sigma$  while  $J_{\theta, \sigma}$  is the portion of the Jacobian that regards  $\theta$  and  $\sigma$ . The estimated value for the variance of  $\sigma$  from the lowest level group (level 2) is used to form the standard error of the residual variance.

However, the variance estimates are not simply the values of  $\theta$  and  $\sigma$  but transformations of that (eq. 15). To estimate the variances of the variance estimates, the Delta method is used so that

$$\text{Var}(\Sigma, \sigma^2) = [\nabla(\Lambda^T \Lambda)]^T \text{Var}(\theta, \sigma^2) [\nabla(\Lambda^T \Lambda)] \quad (116)$$

where the gradient ( $\nabla(\cdot)$ ) is taken with respect to the elements of  $\Sigma$  and  $\sigma^2$ , and  $\text{Var}(\theta, \sigma^2)$  is from eq. 115.

## Model Evaluation: Wald Test

{model-evalua}

We can use the a Wald test to test both fixed effects parameters ( $\beta$ ) and variance of the random parameters ( $\Lambda$ ).

The Wald test compares estimated parameters with null hypothesis values. In the default case the null hypothesis is that value of the parameters is 0.

In this default case, if the test fails to reject the null hypothesis, removing the variables from the model will not substantially harm the fit of that model.

One advantage of the Wald test is that it can be used to test multiple hypotheses about multiple parameters simultaneously.

To test  $q$  hypotheses on  $p$  estimated parameters, let  $\hat{P}$  be the vector of estimated coefficients,  $R$  be a  $q \times p$  hypothesis matrix (this matrix has 1 row per coefficient being tested with a value of 1 in the column

corresponding to that coefficient),  $\hat{V}$  be the estimated covariance matrix for  $\hat{P}$ , and  $r$  be the vector of hypothesized values for  $\hat{\beta}$ .

Then the Wald test statistic for multiple parameters is equal to:

$$W = (R\hat{\beta} - r)'(R\hat{V}R')^{-1}(R\hat{\beta} - r)$$

The resulting test statistic can be tested against a chi-square distribution. For this test, the degrees of freedom is the number of parameters that are tested.

$$W \sim \chi^2(p)$$

## References

{references}

Bates, D., Machler, M., Bolker, B. M., & Walker, S. C. (2015), Fitting linear mixed-effects models using lme4. *Journal of Statistical Software*, 67(1), 1–48.

Bates, D., & Pinheiro, J. C. (1998). *Computational methods for multilevel modelling*. Bell Labs Report.

Binder, D. A. (1983). On the variances of asymptotically normal estimators from complex surveys. *International Statistical Review*, 51(3), 279–292.

Gelman, A. (2007). Struggles with survey weighting and regression modeling. *Statistical Science*, 22(2), 153–164.

Henderson, C. R. (1982). Analysis of covariance in the mixed model: higher-level, nonhomogeneous, and random regressions. *Biometrics*, 38(3), 623–640.

Integration by Substitution. (n.d.). In Wikipedia. Retrieved February 13, 2019, from [https://en.wikipedia.org/wiki/Integration\\_by\\_substitution](https://en.wikipedia.org/wiki/Integration_by_substitution)

Rabe-Hesketh, S., & Skrondal, A. (2006). Multilevel modelling of complex survey data. *Journal of the Royal Statistical Society. Series A (Statistics in Society)*, 169(4), 805–827.

Rabe-Hesketh, S., Skrondal, A., & Pickles, A. (2002). Reliable estimation of generalized linear mixed models using adaptive quadrature. *Stata Journal*, 2, 1–21.

Trefethen, L. N., & Bau, D. (1997). *Numerical linear algebra*. Philadelphia, PA: SIAM.