

Nonparametric Estimation in Linear Mixed Models with Uncorrelated Homoscedastic Errors

Eugène-Patrice Ndong Nguéma, Bertrand Fesuh Nono*, Henri Gwét

Department of Mathematics, Ecole Polytechnique, Yaoundé, Cameroon

Email: *fesuhbe@yahoo.co.uk

How to cite this paper: Ndong Nguéma, E.-P., Fesuh Nono, B. and Gwét, H. (2021) Nonparametric Estimation in Linear Mixed Models with Uncorrelated Homoscedastic Errors. *Open Journal of Statistics*, 11, 558-605. <https://doi.org/10.4236/ojs.2021.114035>

Received: July 2, 2021

Accepted: August 27, 2021

Published: August 30, 2021

Copyright © 2021 by author(s) and Scientific Research Publishing Inc. This work is licensed under the Creative Commons Attribution International License (CC BY 4.0).

<http://creativecommons.org/licenses/by/4.0/>



Open Access

Abstract

Today, Linear Mixed Models (LMMs) are fitted, mostly, by assuming that random effects and errors have Gaussian distributions, therefore using Maximum Likelihood (ML) or REML estimation. However, for many data sets, that double assumption is unlikely to hold, particularly for the random effects, a crucial component in which assessment of magnitude is key in such modeling. Alternative fitting methods not relying on that assumption (as ANOVA ones and Rao's MINQUE) apply, quite often, only to the very constrained class of variance components models. In this paper, a new computationally feasible estimation methodology is designed, first for the widely used class of 2-level (or longitudinal) LMMs with only assumption (beyond the usual basic ones) that residual errors are uncorrelated and homoscedastic, with no distributional assumption imposed on the random effects. A major asset of this new approach is that it yields nonnegative variance estimates and covariance matrices estimates which are symmetric and, at least, positive semi-definite. Furthermore, it is shown that when the LMM is, indeed, Gaussian, this new methodology differs from ML just through a slight variation in the denominator of the residual variance estimate. The new methodology actually generalizes to LMMs a well known nonparametric fitting procedure for standard Linear Models. Finally, the methodology is also extended to ANOVA LMMs, generalizing an old method by Henderson for ML estimation in such models under normality.

Keywords

Clustered Data, Linear Mixed Model, Fixed Effect, Uncorrelated Homoscedastic Error, Random Effects Predictor

1. Introduction

The Linear Mixed Model (LMM) is an extension of the classical Linear Model

(LM) aimed at modeling a continuous scalar response Y in terms of observed covariates some of which (say X_1, \dots, X_p) have fixed effects as in the LM, being the same for all individuals in the population under study, and others (say Z_1, \dots, Z_r) have random effects, thus possibly varying between some well identified subgroups in that population. Mixed models have been used to analyze data sets with clustered, longitudinal or multilevel structure in a variety of fields, such as medicine [1] [2], agriculture [3], animal breeding [4], small area estimation [5] [6], genetics [7] [8], growth modeling [9] [10], etc. Detailed presentations of LMMs can be found in [11] or Demidenko [12], with a more practical emphasis in West *et al.* [13]. The subclass of variance components (or ANOVA) LMMs is thoroughly examined in Searle *et al.* [14].

At the outset of this field, that latter class of LMMs was almost exclusively the only one considered to fit real world data sets, with ANOVA methods at the forefront. But things started to turn around by the end of the 1970s, with the advent of widespread powerful computational tools (both in hardware and software) which have now sufficiently matured and are widely accessible. As a consequence, nowadays by far the most popular approach to fitting an LMM to a given data set is to use a Gaussian LMM, *i.e.* assuming that both the random effects and the model residual errors have Gaussian distributions, then using either the Maximum Likelihood (ML) or the Restricted ML (REML) estimation procedures. This is not easy task, however, since the corresponding (restricted) likelihood equations are quite involved, even to simply be derived, and have no closed form solutions. Thus in the last 4 - 5 decades, considerable endeavor, both theoretical and computational, has been devoted to deriving and numerically solving these equations as efficiently as possible [11] [15] [16] [17] [18] [19]. The numerical solution of the nonlinear system of (restricted) likelihood equations is done, most often, through using one of two standard iterative methods: the Newton-Raphson (NR) or Fisher scoring (FS, a more statistical variant of NR) and the Expectation-Maximization (EM) algorithm in various forms [20] [21] [22].

Despite that unrivaled popularity, there is an obvious major issue with fitting Gaussian LMMs: the assumption of normality of both random effects and errors is dubious in many practical settings, especially for the former. Lange and Ryan [23], for instance, present practical cases of nonnormal random effects. In the last decade, faced with bigger and bigger data sets, both in size and dimensionality, a growing interest has been focused on how to analyze LMMs for such data. For instance, a quasi-likelihood approach for estimation and inference in linear mixed-effects models with high-dimensional fixed effects and possibly large or unbalanced cluster sizes has been recently proposed by Li *et al.* [24].

Now, prior to the widespread adoption of ML/REML in LMMs fitting, some methods not relying on Gaussian assumptions had been developed and used such as ANOVA methods [14] [25] [26] [27] [28], Henderson's Methods I, II and III [29], Rao's Minimum Norm Quadratic Unbiased Estimation (MINQUE) method [30] [31] [32] [33], iterative MINQUE [34]. But these are tailored only

to a very specific and limited class of LMMs, namely ANOVA ones (or variance components LMMs). Moreover, these methods share the common significant drawback of not ensuring nonnegative estimates of variance components. A qualified discussion in Searle *et al.* [14] seldom recommended them for general usage. Alternatively, a quasi-likelihood approach is presented by Jiang [35] using the REML equations in non-Gaussian ANOVA LMMs, showing consistency and asymptotic normality of the variance components estimators under some additional assumptions. Heyde [36] does the same, but emphasizing the relationship between quasi-likelihood and estimating functions. Iterative weighted least squares and iterative generalized least squares estimation methods for LMMs are presented, respectively, in Jiang *et al.* [37] and Goldstein [38].

Being based on quite restrictive assumptions, the aforementioned methods developed for non-Gaussian LMMs have a limited range of LMMs for which they are theoretically valid fitting methods. For instance, the overly simplifying assumptions on covariance matrices structure (like the ones in ANOVA and MINQUE methods) even exclude most Gaussian LMMs from that validity range, which helps explain why the ML and REML approaches have superseded them. In contrast, our goal in this work is to design an estimation methodology for LMMs in which the fixed effects vector of parameters and the random effects covariance matrix are estimated based on assumptions as weak as possible. We will first restrict attention to the very popular subclass of 2-level (also called *longitudinal*) LMMs pioneered by Laird and Ware [22], before extending the approach to fit ANOVA LMMs as well.

To achieve that goal, we take the practical standpoint here that, in most situations, except when some additional information is available about the data (such as serial correlation in errors for some longitudinal data), one has no other choice than to assume that the residual errors in the LMM to fit are uncorrelated and homoscedastic (u.h.o.). This is probably the default option for error modeling (with the Gaussian assumptions when using ML or REML) in most statistical software packages tailored for parameters estimation in an LMM. In our modeling of 2-level LMMs here, based on adding only that assumption to the basic ones of zero mean and finite covariances for random effects and errors, we devise a new approach for estimating the fixed effects parameters vector β , the cluster random effects covariance matrix \mathbf{D} and the residual errors variance σ_ϵ^2 . Thus, we do not impose any assumption on the type of clusters random effects distribution. The approach can be viewed as an adaptation to 2-level LMMs of the well known 2-step procedure for fitting the LM to a given data set, under that same assumption, where first β is estimated by Ordinary Least Squares (OLS), then the variance σ_ϵ^2 of the errors by a carefully designed unbiased estimator. On closer scrutiny, it also turns out to generalize a little popularized estimating procedure for ML computation in Gaussian variance components models credited to Henderson by Harville [16] and detailed in Searle *et al.* ([14], pages 278-279).

As for the organization of the material presented in this article, Section 2 briefly reviews LMMs, up to the random effects prediction issue. Section 3 highlights some key preliminary results on LMMs which do not involve any Gaussian distributional assumptions, which will serve as backbone for our new estimation methodology to be designed. In Section 4, we construct that new estimation methodology for 2-level LMMs with u.h.o. errors. Section 5 shows that our estimation methodology can be adapted to both 2-level LMMs with u.h.o. errors and a diagonal cluster random effects covariance matrix and ANOVA LMMs. In Section 6, our new estimation approach is compared with the traditional Gaussian ML through a simulation study, an application to a classical data set with cluster structure, and a longitudinal data, with implementations done in the R software [39]. Section 7 draws our conclusion about the work presented. The Appendix contains the most lengthy proof of a result presented in the running text. A supplementary material document is also provided, which gathers some known results on LMMs which we have used, but scattered here and there in the literature, and some new results of our own, as well as important implementation details about our presented iterative methods for fitting LMMs. To make a distinction with the article, the numbering of sections, results and equations in that document is prefixed by the letter “S”.

Before we proceed, please note that in this paper, all vectors are columns. Moreover, A^T , $\text{tr}(A)$, $\|u\| = \sqrt{u^T u}$ and $\|u\|_M = \sqrt{u^T M u}$ respectively denote transpose of matrix A , trace of square matrix A , Euclidean norm and norm w.r.t. the SPD matrix M of vector u , where “SPD” stands for *symmetric and positive definite* while “SPSD” means *symmetric and positive semi-definite* (matrix). \mathbf{I}_n is the identity matrix of order n . $\mathcal{M}_{n,p}(\mathbb{R})$ and $\mathcal{M}_n(\mathbb{R})$ are the spaces of matrices with real elements, respectively $n \times p$ and $n \times n$. $\mathbb{E}(X)$ and $\mathbb{M}\text{cov}(X)$ denote the mean and covariance matrix of random vector (or variable) X . $\mathcal{N}_q(\mathbf{m}, \Sigma)$ is the q -dimensional Gaussian distribution with mean vector \mathbf{m} and covariance matrix Σ , $\stackrel{\mathcal{L}}{\sim}$ is for “has probability distribution”, while *estimating equation* will be abbreviated *EE*.

2. Linear Mixed Models: An Overview

2.1. The General Form of an LMM

The general form of an LMM is [11]:

$$\mathbf{Y} = \mathbf{X}\beta + \mathbf{Z}U + \varepsilon, \quad (1)$$

where $\mathbf{Y} \in \mathbb{R}^n$ is the observed response, $\beta \in \mathbb{R}^p$ is the unknown vector of fixed effects parameters, $U \in \mathbb{R}^q$ is the vector of unobserved random effects, $\varepsilon \in \mathbb{R}^n$ comprises the unknown residual errors, while $\mathbf{X} \in \mathcal{M}_{n,p}(\mathbb{R})$ and $\mathbf{Z} \in \mathcal{M}_{n,q}(\mathbb{R})$ are given design matrices. The usual basic assumptions are [11] [12] [22]:

Assumption $\mathcal{A}1_g$. U and ε are two independent zero mean random vectors with respective positive definite covariance matrices $\mathbb{M}\text{cov}(U) = \mathbf{G} \in \mathcal{M}_q(\mathbb{R})$

and $\text{M cov}(\varepsilon) = \mathbf{R} \in \mathcal{M}_n(\mathbb{R})$.

With these assumptions, the covariance matrix of \mathbf{Y} is $\mathbf{V} = \text{M cov}(\mathbf{Y}) = \mathbf{ZGZ}^T + \mathbf{R}$. We will also use the following assumption:

Assumption $\mathcal{A}2_g$. The $n \times p$ design matrix \mathbf{X} in (1) has full column rank.

In (1), the main objective is to estimate β and the covariance matrices \mathbf{R} and \mathbf{G} . LMMs are classified into two main groups, based on Gaussian distributional assumptions made or not on the random effects and errors in devising estimation procedures. Model (1) is termed Gaussian if, in addition to Assumption $\mathcal{A}1_g$, it holds that:

Assumption $\mathcal{A}3_g$. $U \stackrel{\mathcal{L}}{\sim} \mathcal{N}_q(\mathbf{0}, \mathbf{G})$ and $\varepsilon \stackrel{\mathcal{L}}{\sim} \mathcal{N}_n(\mathbf{0}, \mathbf{R})$.

If the random effects and/or errors are not assumed to be normal, but Assumption $\mathcal{A}1_g$ holds, then model (1) is usually said to be (rather abusively) a non-Gaussian LMM. The main goal in this work is to devise estimation procedures for LMMs which do not use Assumption $\mathcal{A}3_g$. The only key added assumption will rather be:

Assumption $\mathcal{A}4_g$. $\mathbf{R} = \sigma_\varepsilon^2 \mathbf{I}_n$ ($\sigma_\varepsilon^2 > 0$), i.e. the residual errors in (1) are uncorrelated and homoscedastic.

2.2. The 2-Level (or Longitudinal) LMM

We will first design our new estimation methodology for one of the types of LMM most used in statistical analyzes. Since Laird and Ware [22], this type of LMM is usually qualified as *longitudinal*, but we feel it more encompassing to label them as 2-level. In this article, we will use the *units-cluster* terminology for 2-level LMMs. At the cluster level, such a model can be written by grouping the scalar responses for all units in each cluster j as:

$$\mathbf{Y}_j = \mathbf{X}_j \beta + \mathbf{Z}_j U_j + \varepsilon_j, \quad j = 1, \dots, m, \tag{2}$$

where $\mathbf{Y}_j = (Y_{1j}, \dots, Y_{n_j, j})^T \in \mathbb{R}^{n_j}$ is the response vector; $\varepsilon_j = (\varepsilon_{1j}, \dots, \varepsilon_{n_j, j})^T \in \mathbb{R}^{n_j}$ is the vector of the errors; $\mathbf{X}_j \in \mathcal{M}_{n_j, p}(\mathbb{R})$ is the fixed effects design matrix; $\mathbf{Z}_j \in \mathcal{M}_{n_j, r}(\mathbb{R})$ is the random effects design matrix. The unknown parts of the model are $\beta = (\beta_1, \beta_2, \dots, \beta_p)^T \in \mathbb{R}^p$, the vector of fixed effects parameters, and $U_j \in \mathbb{R}^r$, the unobserved vector of random effects for cluster j . We will also consider the following usual basic assumptions for such models:

Assumption $\mathcal{A}1$. The U_j 's are independent and identically distributed in \mathbb{R}^r .

Assumption $\mathcal{A}2$. The ε_j 's are independent.

Assumption $\mathcal{A}3$. The set of U_j 's and set of ε_j 's are independent from each other.

Assumption $\mathcal{A}4$. $\mathbb{E}(U_j) = \mathbf{0} \in \mathbb{R}^r$, $\mathbb{E}(\varepsilon_j) = \mathbf{0} \in \mathbb{R}^{n_j}$, $\text{M cov}(U_j) = \mathbf{D} \in \mathcal{M}_r(\mathbb{R})$, $\text{M cov}(\varepsilon_j) = \mathbf{R}_j \in \mathcal{M}_{n_j}(\mathbb{R})$. Here, \mathbf{D} and \mathbf{R}_j are, respectively, the covariance matrices of U_j and ε_j and are assumed finite and positive definite.

Model (2) can be written in the form (1), but with the covariance matrices having special diagonal block structures: $\mathbf{G} = \text{diag}(\mathbf{D}, \dots, \mathbf{D})$, $\mathbf{R} = \text{diag}(\mathbf{R}_1, \dots, \mathbf{R}_m)$, $\mathbf{V} = \text{diag}(\mathbf{V}_1, \dots, \mathbf{V}_m)$, with $\mathbf{V}_j = \mathbf{Z}_j \mathbf{D} \mathbf{Z}_j^T + \mathbf{R}_j$, $j = 1, \dots, m$.

2.3. Variance Components Models (or ANOVA LMMs)

We will subsequently show that our methodology can also be adapted to the popular ANOVA LMMs, *i.e.* LMMs (1) in which the term \mathbf{ZU} can be decomposed as:

$$\mathbf{ZU} = \mathbf{Z}_1 U_1 + \dots + \mathbf{Z}_m U_m, \quad (3)$$

where $\mathbf{Z}_1, \dots, \mathbf{Z}_m$ are given design matrices, with $\mathbf{Z}_j \in \mathcal{M}_{n, q_j}(\mathbb{R})$, $j = 1, \dots, m$, and U_1, \dots, U_m now become the random effects vectors of respective dimensions q_1, \dots, q_m ($q_1 + \dots + q_m = q$). Taking $\mathbf{Z}_{m+1} = \mathbf{I}_n$ and $U_{m+1} = \varepsilon$ (so $q_{m+1} = n$), the ANOVA LMM (1) with random effects term (3) is more neatly written:

$$\mathbf{Y} = \mathbf{X}\beta + \mathbf{Z}\tilde{U}, \quad (4)$$

with $\mathbf{Z} = (\mathbf{Z}_1, \dots, \mathbf{Z}_{m+1})$, $\tilde{U} = (U_1^T, \dots, U_{m+1}^T)^T$ and default assumptions:

M1. U_1, \dots, U_{m+1} are independent;

M2. $\forall j = 1, \dots, m+1$, $\mathbb{E}(U_j) = \mathbf{0}$ and $\mathbb{M}\text{cov}(U_j) = \sigma_j^2 \mathbf{I}_{q_j}$, with $\sigma_j^2 > 0$ finite.

Then (4) implies $\mathbf{G} = \text{diag}(\sigma_1^2 \mathbf{I}_{q_1}, \dots, \sigma_m^2 \mathbf{I}_{q_m})$, and the only unknown parameters in (4), besides β , are the vector of components variances $\theta = (\sigma_1^2, \dots, \sigma_{m+1}^2)^T$. That is why (4) is called, under **M1** and **M2**, a *variance components model*.

2.4. Prediction of Random Effects: The BLP and the BLUP of U

Besides being sometimes an important target in LMM modeling, it turns out that having an effective and computable (given parameters values) predictor for the random effects vector U in the LMM (1) is a crucial component in our newly designed estimation methodology. In that respect, a first candidate is the best linear predictor (BLP) of U given \mathbf{Y} (see [14], Chapter 7):

$$\text{BLP}(U | \mathbf{Y}) = \bar{U} = \mathbf{GZ}^T \mathbf{V}^{-1} (\mathbf{Y} - \mathbf{X}\beta). \quad (5)$$

But a far more popular predictor is the best linear unbiased predictor (BLUP) of $U | \mathbf{Y}$ given, under Assumptions $\mathcal{A}1_g - \mathcal{A}2_g$, by Searle [40]:

$$\text{BLUP}(U | \mathbf{Y}) = \tilde{U}^* = \mathbf{GZ}^T \mathbf{V}^{-1} (\mathbf{Y} - \mathbf{X}\tilde{\beta}^*), \quad (6)$$

where $\tilde{\beta}^*$ is the Best Linear Unbiased Estimator (BLUE) of β ,

$$\tilde{\beta}^* = (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1} \mathbf{Y}. \quad (7)$$

But, in our new methodology (as in previous ones), trying to use either the BLP or the BLUP requires to first estimate \mathbf{G} and \mathbf{R} . Replacing them in (6) by estimators $\hat{\mathbf{G}}$ and $\hat{\mathbf{R}}$ yields $\hat{U}^* = \text{EBLUP}(U | \mathbf{Y}) = \hat{\mathbf{G}} \mathbf{Z}^T \hat{\mathbf{V}}^{-1} (\mathbf{Y} - \mathbf{X}\hat{\beta}^*)$, a so called Empirical BLUP (EBLUP) of $U | \mathbf{Y}$, with $\hat{\mathbf{V}} = \mathbf{Z} \hat{\mathbf{G}} \mathbf{Z}^T + \hat{\mathbf{R}}$, and an Empirical BLUE (EBLUE) $\hat{\beta}^* = (\mathbf{X}^T \hat{\mathbf{V}}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \hat{\mathbf{V}}^{-1} \mathbf{Y}$ of β . However, importantly for designing our new fitting methodology for LMMs, we stress that

the $\text{BLUP}(U | \mathbf{Y})$ can also be viewed as a preliminary estimator of $\text{BLP}(U | \mathbf{Y})$. Consequently, any $\text{EBLUP}(U | \mathbf{Y})$ can also be taken as a practical estimator of $\bar{U} = \text{BLP}(U | \mathbf{Y})$, thus qualifying also to be an Empirical BLP (EBLP) of $U | \mathbf{Y}$. Indeed, that is how we will estimate \bar{U} , using Henderson's mixed model equations briefly reviewed next.

2.5. Henderson's Mixed Model Equations

For given \mathbf{G} and \mathbf{R} in a Gaussian LMM (1), to simultaneously get an estimator for β and a predictor for U , Henderson [41] [42] maximized, w.r.t. β and U (the latter also viewed as a parameter), the joint density of the random couple (\mathbf{Y}, U) . That maximization yields the so called *Henderson's mixed model equations* (HMMEs) given in matrix form as:

$$\begin{pmatrix} \mathbf{X}^T \mathbf{R}^{-1} \mathbf{X} & \mathbf{X}^T \mathbf{R}^{-1} \mathbf{Z} \\ \mathbf{Z}^T \mathbf{R}^{-1} \mathbf{X} & \mathbf{Z}^T \mathbf{R}^{-1} \mathbf{Z} + \mathbf{G}^{-1} \end{pmatrix} \begin{pmatrix} \tilde{\beta} \\ \tilde{U} \end{pmatrix} = \begin{pmatrix} \mathbf{X}^T \mathbf{R}^{-1} \mathbf{Y} \\ \mathbf{Z}^T \mathbf{R}^{-1} \mathbf{Y} \end{pmatrix}. \quad (8)$$

From the outset, there has been debate over whether this was a valid way of trying to estimate β and predict U . Henderson himself acknowledged (see [43], page 16) that (even given an observed $\mathbf{Y} = \mathbf{y} \in \mathbb{R}^n$) the maximized function is not a true likelihood for the couple (β, U) because U is not a fixed unknown parameter but rather an unobserved random variable. So, strictly speaking, this does not qualify as an ML method. Nonetheless, from our standpoint, that debate is mostly peripheral. The only two relevant practical questions are: 1) What are the solutions $\tilde{\beta}$ and \tilde{U} of the linear system (8)? 2) What properties do they have if taken as respective estimator of β and predictor of U ? The answer to both questions is provided by the following result, credited by Harville [16] to Henderson in an unpublished 1963 report:

Theorem 1 ([16] [44]). *If \mathbf{X} is full column rank while \mathbf{G} and \mathbf{R} are SPD matrices, then the solutions of the linear system (8) are $\tilde{\beta} = \tilde{\beta}^*$ and $\tilde{U} = \tilde{U}^*$ given by (6)-(7), with $\mathbf{V} = \mathbf{Z}\mathbf{G}\mathbf{Z}^T + \mathbf{R}$.*

So, in the general LMM (1), given \mathbf{G} , \mathbf{R} and the observed response vector \mathbf{Y} , the BLUE $\tilde{\beta}^*$ of β and the BLUP \tilde{U}^* of $U | \mathbf{Y}$ are the unique solutions of the HMMEs.

3. More about Henderson's Mixed Model Equations and LMMs without Gaussian Assumptions

As announced in the introduction, we intend to devise a new estimation methodology, first for 2-level LMMs, which, contrary to usage outside of the Gaussian case, simultaneously estimates unknown parameters and predicts random effects values under Assumptions $\mathcal{A}1$ - $\mathcal{A}4$ and $\mathcal{A}4_g$. In that perspective, we give, in this section, and, especially, in Section 3.2, a series of new results about LMMs which do not use any Gaussian assumption, be it on the random effects or the residual errors. Those results will serve as basis for the design of our new fitting methodology for LMMs in Section 4.

3.1. An Important Preliminary: The BLUE, the BLUP and the HMMEs Do Not Need Gaussian Assumptions

Though Henderson, as reported above, famously discovered his eponymous mixed model Equations (8) coming from the Gaussian route, Theorem 1, which solves them, uses no Gaussian distributional assumptions for that. It actually uses no probabilistic, nor statistical framework or reasoning, whatsoever. It is rather a pure Matrix Algebra result asserting that the linear system (8) has as unique solutions the BLUE $\tilde{\beta}^*$ (7) of β and the BLUP \tilde{U}^* (6) of $U | Y$. Now, one can derive both $\tilde{\beta}^*$ and \tilde{U}^* by optimizing, for each, a specific minimum variance criterion, but without relying on any Gaussian assumption or any parametric one for that matter [40]. That double observation makes it that the HMMEs are not attached to Gaussian distributional assumptions, contrary to what people always feel compelled to set before they use them.

Remark 1. *For the latter reason, we will use the HMMEs as our first two estimating equations when devising the first approach (coded 3S-A1-V1) of our new estimation methodology for LMMs in Section 4.1.*

3.2. More about the HMMEs Solutions

Here, we present a series of new Matrix Algebra results and their consequences for LMMs with u.h.o. errors, without imposing any Gaussian assumption, much in line with Theorem 1. The discovery of those results triggered the design of our new fitting methodology for that class of LMMs. They derive from the very peculiar structure of the first HMME which we start by stressing.

3.2.1. The First HMME as a Weighted Least Squares Problem

In (8), the first equation can be rewritten: $X^T R^{-1} X = X^T R^{-1} (Y - Z\tilde{U})$. One then recognises the normal equations of the Weighted Least Squares (WLS) problem:

$$\min_{\beta \in \mathbb{R}^p} \left\| (Y - Z\tilde{U}) - X\beta \right\|_{R^{-1}}. \quad (9)$$

Combining that observation with Theorem 1 yields the remarkable result:

Theorem 2. *In the LMM (1) with Assumptions \mathcal{A}_{1_g} - \mathcal{A}_{2_g} , and given the BLUP \tilde{U}^* of $U | Y$, the BLUE $\tilde{\beta}^*$ of β is the WLS estimate of β in the LM: $Y - Z\tilde{U}^* = X\beta + \tilde{\varepsilon}$, with response $Y - Z\tilde{U}^*$, design matrix X , error $\tilde{\varepsilon}$ and weighting matrix R^{-1} .*

But Theorem 2 will be more useful to us in designing our new methodology for LMM fitting in Section 4 when one adds the assumption that the LMM has uncorrelated and homoscedastic errors. We detail hereafter why.

3.2.2. More about the BLUE of β and the BLUP of U in an LMM with u.h.o. Errors

Here, we focus attention on LMMs with u.h.o. residual errors, *i.e.* Assumption \mathcal{A}_{4_g} holds, in addition to \mathcal{A}_{1_g} - \mathcal{A}_{2_g} . But again, no Gaussian assumption will be relied upon: only the u.h.o. errors and their impact on the geometry of the

HMMEs solutions play a role here. Furthermore, we emphasize that our results to be presented hereafter are valid for any LMM (1) with u.h.o. errors, and not only for 2-level LMMs with such errors. In that respect, first note that in the u.h.o. residual errors scenario, the HMMEs (8) simplify significantly. Indeed, if $\mathbf{R} = \sigma_\varepsilon^2 \mathbf{I}_n$ with $\sigma_\varepsilon^2 > 0$, then (8) reduces to:

$$\begin{pmatrix} \mathbf{X}^T \mathbf{X} & \mathbf{X}^T \mathbf{Z} \\ \mathbf{Z}^T \mathbf{X} & \mathbf{Z}^T \mathbf{Z} + \sigma_\varepsilon^2 \mathbf{G}^{-1} \end{pmatrix} \begin{pmatrix} \tilde{\beta} \\ \tilde{U} \end{pmatrix} = \begin{pmatrix} \mathbf{X}^T \mathbf{Y} \\ \mathbf{Z}^T \mathbf{Y} \end{pmatrix}. \tag{10}$$

The results to be derived hereafter highlight some peculiar properties about the solutions of the linear system (10), thus yielding some remarkable intertwined relationships, under Assumption $\mathcal{A}4_g$, between its unique solutions (Theorem 1), the BLUE $\tilde{\beta} = \tilde{\beta}^*$ of β and $\tilde{U} = \tilde{U}^*$, the BLUP of $U | Y$.

We start with an identity which is a trivial consequence of the structure of V in an LMM (1) under u.h.o. errors. Indeed, then $V = \mathbf{ZGZ}^T + \sigma_\varepsilon^2 \mathbf{I}_n$, so:

$$\mathbf{I}_n - \mathbf{ZGZ}^T V^{-1} = \sigma_\varepsilon^2 V^{-1}. \tag{11}$$

This is instrumental first for:

Lemma 3. Consider the LMM (1) with Assumptions $\mathcal{A}1_g$ and $\mathcal{A}4_g$. Let $(\tilde{\beta}, \tilde{U}) \in \mathbb{R}^p \times \mathbb{R}^q$ such that $\tilde{U} = \mathbf{GZ}^T V^{-1} (\mathbf{Y} - \mathbf{X}\tilde{\beta})$. Then one has:

$$\mathbf{Y} - \mathbf{X}\tilde{\beta} - \mathbf{Z}\tilde{U} = \sigma_\varepsilon^2 V^{-1} (\mathbf{Y} - \mathbf{X}\tilde{\beta}). \tag{12a}$$

Moreover, if Assumption $\mathcal{A}2_g$ also holds, then the following equivalence is true:

$$\tilde{\beta} = (\mathbf{X}^T V^{-1} \mathbf{X})^{-1} \mathbf{X}^T V^{-1} \mathbf{Y} \Leftrightarrow \tilde{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{Y} - \mathbf{Z}\tilde{U}). \tag{12b}$$

Proof. Let Assumptions $\mathcal{A}1_g$ and $\mathcal{A}4_g$ hold, $\mathbf{R} = \sigma_\varepsilon^2 \mathbf{I}_n$ and $\tilde{U} = \mathbf{GZ}^T V^{-1} \tilde{Y}$, with $\tilde{Y} = \mathbf{Y} - \mathbf{X}\tilde{\beta}$, for an arbitrary $\tilde{\beta} \in \mathbb{R}^p$. Then $\mathbf{Y} - \mathbf{X}\tilde{\beta} - \mathbf{Z}\tilde{U} = \tilde{Y} - \mathbf{ZGZ}^T V^{-1} \tilde{Y} = (\mathbf{I}_n - \mathbf{ZGZ}^T V^{-1}) \tilde{Y} = \sigma_\varepsilon^2 V^{-1} \tilde{Y}$, the last equality thanks to (11), which proves (12a).

Now, if Assumption $\mathcal{A}2_g$ is also true, then the square matrix $\mathbf{X}^T \mathbf{X}$ is nonsingular; so, given that $\mathbf{Y} - \mathbf{Z}\tilde{U} = \mathbf{X}\tilde{\beta} + \sigma_\varepsilon^2 V^{-1} \tilde{Y}$ from (12a), then $(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{Y} - \mathbf{Z}\tilde{U}) = \tilde{\beta} + \sigma_\varepsilon^2 (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T V^{-1} (\mathbf{Y} - \mathbf{X}\tilde{\beta})$. Hence, since $\sigma_\varepsilon^2 > 0$, $\tilde{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{Y} - \mathbf{Z}\tilde{U}) \Leftrightarrow \mathbf{X}^T V^{-1} (\mathbf{Y} - \mathbf{X}\tilde{\beta}) = \mathbf{0}$. Then (12b) follows because Assumption $\mathcal{A}2_g$ implies that $\mathbf{X}^T V^{-1} \mathbf{X}$ is an SPD matrix, thus nonsingular. \square

Theorem 1 and Lemma 3 imply the following for the solutions of (10):

Theorem 4. In the LMM (1), if $P = V^{-1} - V^{-1} \mathbf{X} (\mathbf{X}^T V^{-1} \mathbf{X})^{-1} \mathbf{X}^T V^{-1}$ and Assumptions $\mathcal{A}1_g - \mathcal{A}2_g$, $\mathcal{A}4_g$ hold, then:

$$\mathbf{Y} - \mathbf{X}\tilde{\beta}^* - \mathbf{Z}\tilde{U}^* = \sigma_\varepsilon^2 V^{-1} (\mathbf{Y} - \mathbf{X}\tilde{\beta}^*) = \sigma_\varepsilon^2 P \mathbf{Y}, \tag{13a}$$

$$\tilde{\beta}^* = (\mathbf{X}^T V^{-1} \mathbf{X})^{-1} \mathbf{X}^T V^{-1} \mathbf{Y} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{Y} - \mathbf{Z}\tilde{U}^*). \tag{13b}$$

Proof. Thanks to Lemma 3, the respective definitions (7) and (6) of the BLUE $\tilde{\beta}^*$ and the BLUP \tilde{U}^* , and the standard identity $V^{-1} (\mathbf{Y} - \mathbf{X}\tilde{\beta}^*) = P \mathbf{Y}$. \square

Remark 2. The usefulness of the second expression of $\tilde{\beta}^*$ in (13b) stems

from the fact that one recognizes the OLS, under Assumption $\mathcal{A}2_g$, of β in the LM $\tilde{Y} = X\beta + \tilde{\varepsilon}$, where $\tilde{Y} = Y - Z\tilde{U}^*$. This clearly suggests how, in an iterative algorithm for LMM fitting, one may update the estimate of β given the current estimate of the BLUP of U . That observation is precisely what motivated the construction of the second approach (coded 3S-A1-V2) in our new estimation methodology for LMMs in Section 4.1.

We then have the following corollary:

Corollary 1. If $A = Y^T(Y - X\tilde{\beta}^* - Z\tilde{U}^*)$ and $B = \|Y - X\tilde{\beta}^* - Z\tilde{U}^*\|_G^2$ in the LMM(1) with Assumptions $\mathcal{A}1_g$ - $\mathcal{A}2_g$ and $\mathcal{A}4_g$, then

$$A - B = \sigma_\varepsilon^2 \cdot \|Z^T V^{-1}(Y - X\tilde{\beta}^*)\|_G^2 = \sigma_\varepsilon^2 \cdot \|\tilde{U}^*\|_{G^{-1}}^2. \tag{14}$$

Therefore, the real random variable $Y^T(Y - X\tilde{\beta}^* - Z\tilde{U}^*)$ is always nonnegative.

Proof. Let Assumptions $\mathcal{A}1_g$ - $\mathcal{A}2_g$ and $\mathcal{A}4_g$ be true and $\tilde{Y} = Y - X\tilde{\beta}^*$. First,

$$A - B = (X\tilde{\beta}^* + Z\tilde{U}^*)^T (\tilde{Y} - Z\tilde{U}^*) = \tilde{U}^{*T} Z^T (\tilde{Y} - Z\tilde{U}^*), \tag{15}$$

where we have used the fact that Theorem 4 and Remark 2 imply:

$X\tilde{\beta}^* \perp \tilde{Y} - Z\tilde{U}^*$. Now, thanks to (13a) in Theorem 4, $\tilde{Y} - Z\tilde{U}^* = \sigma_\varepsilon^2 V^{-1} \tilde{Y}$. The latter, inserted in (15) alongside the expression (6) of \tilde{U}^* , implies:

$$A - B = \sigma_\varepsilon^2 \tilde{Y}^T V^{-1} Z G Z^T V^{-1} \tilde{Y} = \sigma_\varepsilon^2 (Z^T V^{-1} \tilde{Y})^T G (Z^T V^{-1} \tilde{Y}) = \sigma_\varepsilon^2 \|Z^T V^{-1} \tilde{Y}\|_G^2.$$

$$\text{Or, } A - B = \sigma_\varepsilon^2 (G Z^T V^{-1} \tilde{Y})^T G^{-1} (G Z^T V^{-1} \tilde{Y}) = \sigma_\varepsilon^2 \tilde{U}^{*T} G^{-1} \tilde{U}^* = \sigma_\varepsilon^2 \|\tilde{U}^*\|_{G^{-1}}^2. \quad \square$$

We end this series of results about the geometry of the HMMs solutions with one about two remarkable expectations related to the u.h.o. residual errors variance σ_ε^2 in the LMM (1):

Theorem 5. In the LMM (1) with Assumptions $\mathcal{A}1_g$ - $\mathcal{A}2_g$ and $\mathcal{A}4_g$, one has, with P as in Theorem 4:

$$\mathbb{E}[Y^T(Y - X\tilde{\beta}^* - Z\tilde{U}^*)] = \sigma_\varepsilon^2 \cdot (n - p), \tag{16a}$$

$$\mathbb{E}\|Y - X\tilde{\beta}^* - Z\tilde{U}^*\|_G^2 = \sigma_\varepsilon^4 \cdot \text{tr}(P). \tag{16b}$$

Proof. First, (13a) in Theorem 4 implies that $Y^T(Y - X\tilde{\beta}^* - Z\tilde{U}^*) = \sigma_\varepsilon^2 Y^T P Y$. Hence, with $m_Y = \mathbb{E}(Y) = X\beta$ and $V = \mathbb{M}\text{cov}(Y)$,

$$\mathbb{E}[Y^T(Y - X\tilde{\beta}^* - Z\tilde{U}^*)] = \sigma_\varepsilon^2 \mathbb{E}(Y^T P Y) = \sigma_\varepsilon^2 [m_Y^T P m_Y + \text{tr}(P V)]. \tag{17a}$$

But $P m_Y = P X \beta = V^{-1} X \beta - V^{-1} X \beta = \mathbf{0}$, so $m_Y^T P m_Y = 0$. Furthermore, since $P V = I_n - V^{-1} X (X^T V^{-1} X)^{-1} X^T$, then $\text{tr}(P V) = n - \text{tr}[(X^T V^{-1} X)^{-1} X^T V^{-1} X] = n - p$. Substituting these results in (17a) yields (16a).

Similarly, from (13a), we get $\|Y - X\tilde{\beta}^* - Z\tilde{U}^*\|_G^2 = \sigma_\varepsilon^4 Y^T P^2 Y$, since P is symmetric. Hence

$$\mathbb{E}\|Y - X\tilde{\beta}^* - Z\tilde{U}^*\|_G^2 = \sigma_\varepsilon^4 \mathbb{E}(Y^T P^2 Y) = \sigma_\varepsilon^4 \cdot [m_Y^T P^2 m_Y + \text{tr}(P^2 V)].$$

Now, $m_y^T P^2 m_y = 0$ as $P m_y = \mathbf{0}$. Also, since $PVP = P$, $\text{tr}(P^2 V) = \text{tr}(PVP) = \text{tr}(P)$, hence (16b). □

Remark 3. Identity (16a) is remarkable and seems familiar. Indeed, it is the analogue, for LMMs with u.h.o. errors, of the identity which yields the unbiased estimator of the residual variance in an LM with that same type of errors. And we will use it in much the same way when deriving our new estimation methodology for LMMs with u.h.o. errors. Indeed, in conjunction with the last sentence of Corollary 1, (16a) easily suggests how to compute, in an LMM fitting algorithm, a nonnegative update to estimate σ_ϵ^2 given the current estimates of β and the BLUP of U .

3.3. More about the BLP of $U | Y$ in an LMM

The BLP $\bar{U} = \text{BLP}(U | Y)$, given by (5), will play a key role in our new estimation methodology for LMMs with u.h.o. errors. That is because, first of all, its derivation does not require any Gaussian distributional assumption. Secondly, its computation and that of its covariance matrix are cheaper than those of the BLUP, especially under this scenario of u.h.o. errors. See the needed formulas in the Supplementary material. In the latter, we also derive new expectations formulas specifically based on the BLP (whereas the ones in the previous section, such as (16a), were based on the BLUP) which can be used to obtain EEs for the residual variance in an LMM.

In our new methodology to be presented from Section 4 onwards, the results in Section 3.2 above and the Supplementary material will provide the tools to derive EEs for the β , σ_ϵ^2 and the $\text{BLP}(U | Y)$. On the other hand, to seek an EE for D , we will start from the fact that given $\bar{U}_j = \text{BLP}(U_j | Y_j)$ and $V_j^* = D - \text{M cov}(U_j)$ ($j = 1, \dots, m$), then, under Assumptions $\mathcal{A}1_g - \mathcal{A}4_g$ and $\mathcal{A}2_g$, we get, as an unbiased estimator of D , the $r \times r$ matrix

$$\tilde{D}_1 = \frac{1}{m} \sum_{j=1}^m (V_j^* + \bar{U}_j \bar{U}_j^T). \tag{18}$$

4. Estimation in 2-Level LMMs with u.h.o. Errors: A New Approach

We first present our approach for estimating the 3 main parameters of interest in a 2-level LMM (2) under Assumptions $\mathcal{A}1 - \mathcal{A}4$, $\mathcal{A}2_g$ and $\mathcal{A}4_g$: $\beta \in \mathbb{R}^p$, $\sigma_\epsilon^2 > 0$ and $D \in \mathcal{M}_r(\mathbb{R})$, alongside obtaining a prediction for $U = (U_1^T, \dots, U_m^T)^T \in \mathbb{R}^q$. For that latter aspect, let, for $j = 1, \dots, m$, $\bar{u}_j = \text{BLP}(U_j | Y_j = y_j) = \sigma_\epsilon^{-2} V_j^* Z_j^T (y_j - X_j \beta)$, the BLP of U_j given response $Y_j = y_j \in \mathbb{R}^{n_j}$ in cluster j . We will predict U through estimating

$$\bar{u} = \text{BLP}(U | Y = y) = (\bar{u}_1^T, \dots, \bar{u}_m^T)^T = \sigma_\epsilon^{-2} V^* Z^T (y - X \beta), \tag{19}$$

the BLP of U given response $Y = y \in \mathbb{R}^n$ for the whole sample.

Our approach is based on iteratively solving appropriate EEs for β , σ_ϵ^2 , D , \bar{u} derived using only sound nonparametric estimation principles. Devising

nonparametric EEs from given data to solve a statistical problem is an old idea in Statistics. The method of moments, whenever applicable, is in that category. In the field of LMMs, one may even say that that's where all started with the ANOVA methods for fitting variance components models. But it has turned especially hard to derive EEs applicable to fit a broad class of LMMs, as unconstrained as possible, but not relying on Gaussian assumptions. Attempts in that direction include the iterative weighted least squares and iterative generalized least squares estimation methods for LMMs presented, respectively, in Jiang *et al.* [37] and Goldstein [38]. But they are still not yet sufficiently general. What we propose hereafter is a 3-step construction to get closer to such a goal.

4.1. Two 3-Step Sequences for Estimation in 2-Level LMMs with u.h.o. Errors

The key point in our approach is the fact that the BLP is computationally cheaper to handle than the BLUP (including at the covariance matrix level) while at the same time, as observed at the end of Section 2.4, an empirical BLUP can also be viewed as an empirical BLP. So, to devise an estimation procedure for β , σ_ε^2 , \mathbf{D} and prediction for U from $\mathbf{Y} = \mathbf{y}$, and given the preliminary results in Section 0, our starting ideas are described in what follows (we provide 2 different versions). In them, when $\hat{\mathbf{D}}$ is an estimate of \mathbf{D} , it is understood that $\hat{\mathbf{G}} = \text{diag}(\hat{\mathbf{D}}, \dots, \hat{\mathbf{D}})$ (m times) is the corresponding estimate of \mathbf{G} .

4.1.1. Starting Ideas: Version 1

- **Step 1: Estimating β and \bar{u} , given estimates of σ_ε^2 and \mathbf{D}**

Given preliminary estimates $\hat{\sigma}_\varepsilon^2$ and $\hat{\mathbf{D}}$ of σ_ε^2 and \mathbf{D} , one can obtain estimates $\hat{\beta}$ of β , and \hat{u} of \bar{u} by solving the system of HMMEs (10) using $\sigma_\varepsilon^2 = \hat{\sigma}_\varepsilon^2$ and $\mathbf{G} = \hat{\mathbf{G}}$.

- **Step 2: Improved estimate of σ_ε^2 , given $\hat{\beta}$, \hat{u} from Step 1**

With $\hat{\beta}$ and \hat{u} obtained as above in Step 1, (16a) in Theorem 5 suggests to take

$$\hat{\sigma}_\varepsilon^2 = \frac{\mathbf{y}^\top (\mathbf{y} - \mathbf{X}\hat{\beta} - \mathbf{Z}\hat{u})}{n - p} \quad (20)$$

as a hopefully improved estimate of σ_ε^2 . Using Corollary 1, the way $\hat{\beta}$ and \hat{u} were obtained in Step 1 implies that $\hat{\sigma}_\varepsilon^2$ is a nonnegative real number.

- **Step 3: Improved estimate of \mathbf{D} given preliminary estimates $\hat{\sigma}_\varepsilon^2$, $\tilde{\mathbf{D}}$, \hat{u} of σ_ε^2 , \mathbf{D} , \bar{u}**

Here, we propose an approach to get an improved estimate $\hat{\mathbf{D}}$ of \mathbf{D} from respective preliminary estimates $\hat{\sigma}_\varepsilon^2$, $\tilde{\mathbf{D}}$, \hat{u} of σ_ε^2 , \mathbf{D} , and \bar{u} in (19). We start with the unbiased estimator $\tilde{\mathbf{D}}_1$ of \mathbf{D} given by (18), but is not computable from the available data since it still depends on the unknown parameters, including \mathbf{D} itself. Nevertheless, with the preliminary estimates $\tilde{\mathbf{D}}$ of \mathbf{D} , and $\hat{\sigma}_\varepsilon^2$ of σ_ε^2 , Equation (S:1.7c) in the Supplementary material suggests estimating \mathbf{V}_j^* by $\hat{\mathbf{V}}_j^* = \hat{\sigma}_\varepsilon^2 (\mathbf{Z}_j^\top \mathbf{Z}_j + \hat{\sigma}_\varepsilon^2 \tilde{\mathbf{D}}^{-1})^{-1}$. Then, for $j = 1, \dots, m$, given the available

estimate $\hat{\mathbf{u}}_j$ of $\bar{\mathbf{u}}_j$, (18) suggests as a hopefully improved computable estimate of \mathbf{D} :

$$\hat{\mathbf{D}} = \frac{1}{m} \sum_{j=1}^m (\hat{\mathbf{V}}_j^* + \hat{\mathbf{u}}_j \hat{\mathbf{u}}_j^T). \tag{21}$$

4.1.2. Starting Ideas: Version 2

Here, Step 3 is the same as in Version 1 above, but Steps 1 and 2 rather go as follows:

- **Step 1: Estimation of $\bar{\mathbf{u}}$ and β , given estimates of β , \mathbf{D} , σ_ε^2**

Assume that we have respective preliminary estimates $\tilde{\beta}$, $\hat{\mathbf{D}}$, $\hat{\sigma}_\varepsilon^2$ of β , \mathbf{D} , σ_ε^2 . Given $\mathbf{Y} = \mathbf{y}$ and (19), we estimate $\bar{\mathbf{u}}$ by

$$\hat{\mathbf{u}} = \hat{\sigma}_\varepsilon^{-2} \hat{\mathbf{V}}^* \mathbf{Z}^T (\mathbf{y} - \mathbf{X} \tilde{\beta}), \tag{22}$$

with $\hat{\mathbf{V}}^* = \hat{\sigma}_\varepsilon^2 (\mathbf{Z}^T \mathbf{Z} + \hat{\sigma}_\varepsilon^2 \hat{\mathbf{G}}^{-1})^{-1}$. From Theorem 4, we hope to get an improved estimate of β through:

$$\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{y} - \mathbf{Z} \hat{\mathbf{u}}). \tag{23}$$

Remark 4. This is not the same as Step 1 of Version 1 because $\tilde{\beta}$ appears on the l.h.s. of (22) instead of $\hat{\beta}$. So here, $\hat{\mathbf{u}}$ and $\hat{\beta}$ are not solutions to the HMMEs given $\mathbf{G} = \hat{\mathbf{G}}$ and $\mathbf{R} = \hat{\sigma}_\varepsilon^2 \mathbf{I}_n$.

Remark 5. In (22), we have $\hat{\mathbf{u}} = (\hat{\mathbf{u}}_1^T, \dots, \hat{\mathbf{u}}_m^T)^T$, where for $j = 1, \dots, m$, $\hat{\mathbf{u}}_j \in \mathbb{R}^r$ is given by:

$$\hat{\mathbf{u}}_j = \hat{\sigma}_\varepsilon^{-2} \hat{\mathbf{V}}_j^* \mathbf{Z}_j^T (\mathbf{y}_j - \mathbf{X}_j \tilde{\beta}), \quad \hat{\mathbf{V}}_j^* = \hat{\sigma}_\varepsilon^2 (\mathbf{Z}_j^T \mathbf{Z}_j + \hat{\sigma}_\varepsilon^2 \hat{\mathbf{D}}^{-1})^{-1}. \tag{24}$$

Remark 6. The vector $\hat{\beta}$ given by (23) is the OLS, under Assumption \mathcal{A}_{2_g} , of β in the LM: $\mathbf{y} - \mathbf{Z} \hat{\mathbf{u}} = \mathbf{X} + \tilde{\varepsilon}$, with $\mathbf{y} - \mathbf{Z} \hat{\mathbf{u}}$ as response and $\tilde{\varepsilon}$ as vector of residual errors.

- **Step 2: Improved estimate of σ_ε^2 , given $\hat{\beta}$, $\hat{\mathbf{u}}$ from Step 1, and preliminary estimates of σ_ε^2 , \mathbf{D}**

With Remark 4 above, there is no guarantee that the numerator of (20) is a nonnegative number given $\hat{\mathbf{u}}$ and $\hat{\beta}$ from Step 1 in this Version 2. To ensure a nonnegative estimate of σ_ε^2 , we combine Corollary 1 and Theorem 5 to obtain, from the available estimates $\tilde{\sigma}_\varepsilon^2$ and $\hat{\mathbf{D}}$ of σ_ε^2 and \mathbf{D} , the hopefully improved estimate of σ_ε^2 :

$$\hat{\sigma}_\varepsilon^2 = \frac{\|\mathbf{y} - \mathbf{X} \hat{\beta} - \mathbf{Z} \hat{\mathbf{u}}\|^2 + \tilde{\sigma}_\varepsilon^2 \|\hat{\mathbf{u}}\|_{\hat{\mathbf{G}}^{-1}}^2}{n - p}. \tag{25}$$

Remark 7. From Corollary 1, Theorem 5 and Propositions S1.4-S1.5, S1.8-S1.9, two alternative somewhat more sophisticated computable nonnegative estimates of σ_ε^2 in this context are:

$$\hat{\sigma}_\varepsilon^2 = \frac{\|\mathbf{y} - \mathbf{X} \hat{\beta} - \mathbf{Z} \hat{\mathbf{u}}\|^2 + q \tilde{\sigma}_\varepsilon^2}{n - p + \text{tr}(\hat{\mathbf{G}}^{-1} \hat{\mathbf{V}}^*)} = \frac{\|\mathbf{y} - \mathbf{X} \hat{\beta} - \mathbf{Z} \hat{\mathbf{u}}\|^2 + q \tilde{\sigma}_\varepsilon^2}{n - p + \text{tr}(\hat{\mathbf{D}}^{-1} \sum_{j=1}^m \hat{\mathbf{V}}_j^*)}, \tag{26a}$$

$$\hat{\sigma}_\varepsilon^2 = \frac{\text{tr}(\hat{\mathbf{V}}^* \mathbf{Z}^T \mathbf{Z})}{\|\hat{\mathbf{u}}\|_{\hat{\mathbf{G}}^{-1}}^2} = \sum_{j=1}^m \text{tr}(\hat{\mathbf{V}}_j^* \mathbf{Z}_j^T \mathbf{Z}_j) / \sum_{j=1}^m \|\hat{\mathbf{u}}_j\|_{\hat{\mathbf{B}}^{-1}}^2. \quad (26b)$$

4.2. Targeted Estimates for β , σ_ε^2 , \mathbf{D} and Prediction of U

4.2.1. The Estimating Equations

The reasoning in Steps 1-2-3 of Version 1 above leads us to seek, given the data vector $\mathbf{Y} = \mathbf{y}$, respective estimates $\hat{\beta}$, $\hat{\sigma}_\varepsilon^2$, $\hat{\mathbf{D}}$, of parameters β , σ_ε^2 , \mathbf{D} , and prediction $\hat{\mathbf{u}}$ of U satisfying the system of equations:

$$\begin{pmatrix} \mathbf{X}^T \mathbf{X} & \mathbf{X}^T \mathbf{Z} \\ \mathbf{Z}^T \mathbf{X} & \mathbf{Z}^T \mathbf{Z} + \hat{\sigma}_\varepsilon^2 \hat{\mathbf{G}}^{-1} \end{pmatrix} \begin{pmatrix} \hat{\beta} \\ \hat{\mathbf{u}} \end{pmatrix} = \begin{pmatrix} \mathbf{X}^T \mathbf{y} \\ \mathbf{Z}^T \mathbf{y} \end{pmatrix}, \quad (27a)$$

$$\hat{\sigma}_\varepsilon^2 = \frac{\mathbf{y}^T (\mathbf{y} - \mathbf{X} \hat{\beta} - \mathbf{Z} \hat{\mathbf{u}})}{n - p}, \quad (27b)$$

$$\hat{\mathbf{D}} = \frac{1}{m} \sum_{j=1}^m (\hat{\mathbf{V}}_j^* + \hat{\mathbf{u}}_j \hat{\mathbf{u}}_j^T). \quad (27c)$$

Steps 1-2-3 of Version 2 suggest the following system, with $\hat{\mathbf{D}}$ as in (27c):

$$\hat{\mathbf{u}} = \hat{\sigma}_\varepsilon^{-2} \hat{\mathbf{V}}^* \mathbf{Z}^T (\mathbf{y} - \mathbf{X} \hat{\beta}), \quad (28a)$$

$$\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{y} - \mathbf{Z} \hat{\mathbf{u}}), \quad (28b)$$

$$\hat{\sigma}_\varepsilon^2 = \frac{\|\mathbf{y} - \mathbf{X} \hat{\beta} - \mathbf{Z} \hat{\mathbf{u}}\|^2 + \hat{\sigma}_\varepsilon^2 \|\hat{\mathbf{u}}\|_{\hat{\mathbf{G}}^{-1}}^2}{n - p}. \quad (28c)$$

We so obtain two distinct sets of fixed points EEs for β , σ_ε^2 , \mathbf{D} and $\bar{\mathbf{u}}$. Our estimation approach for the 2-level LMM (2) under Assumptions $\mathcal{A}1$ - $\mathcal{A}4$ and $\mathcal{A}2_g$, $\mathcal{A}4_g$ consists in solving either of them through successive approximations iterative procedure. However, we stress that, from Theorem 4 and Corollary 1, if $\hat{\mathbf{D}}$ is SPD, $\hat{\sigma}_\varepsilon^2 > 0$ and Assumption $\mathcal{A}2_g$ holds, then one has the following logical equivalence:

$$(27a) - (27b) \Leftrightarrow (28a) - (28b) - (28c).$$

So the above two systems of equations are actually equivalent, and, thus, have the same solutions if any. However, they do suggest two different successive approximations algorithms to try to calculate these 4 unknowns.

4.2.2. 3S: A Code Name for the LMMs New Estimation Methodology

Our proposed two estimating algorithms are presented in the next section. Since each of them is based on a 3-step sequence construction, we shall use the code name 3S for this new methodology for fitting LMMs. In the future, we will present other variants of this approach. So the current one is coded 3S-A1, and the two algorithms we present hereafter are coded 3S-A1-V1 and 3S-A1-V2.

4.2.3. Two Iterative Estimating Procedures for 2-Level LMMs with u.h.o. Errors

Given the response vector $\mathbf{y} = (\mathbf{y}_1^T, \dots, \mathbf{y}_m^T)^T \in \mathbb{R}^n$, and the EEs (27a)-(27c), our

first estimating algorithm for a 2-level LMM with u.h.o. errors is:

Algorithm 3S-A1-V1. Estimating β , σ_ε^2 , D and predicting U_1, \dots, U_m in a 2-level LMM with u.h.o. errors. Version 1

1) **Initialization:** At iteration 0, we estimate β , σ_ε^2 , D and predict U_1, \dots, U_m , as follows:

a) $\hat{\beta}^{(0)} = (X^T X)^{-1} X^T y$, the OLS estimate of β in the linear model $y = X\beta + \varepsilon$;

b) $\hat{\sigma}_\varepsilon^{2(0)} = \|y - X\hat{\beta}^{(0)}\|^2 / (n - p)$, the corresponding unbiased estimate of the residual variance;

c) $\forall j = 1, \dots, m$, $\hat{u}_j^{(0)} = (Z_j^T Z_j)^+ Z_j^T y_j^{*(0)}$, the OLS estimate of U_j in the linear model $y_j^{*(0)} = Z_j U_j + \varepsilon_j$, with U_j considered as fixed parameter, Z_j the design matrix, $y_j^{*(0)} = y_j - X_j \hat{\beta}^{(0)}$ as the response, and $(Z_j^T Z_j)^+$ being the Moore-Penrose pseudo-inverse of $Z_j^T Z_j$;

d) $\hat{D}^{(0)} = \frac{1}{m} \sum_{j=1}^m \hat{u}_j^{(0)} \hat{u}_j^{(0)T}$;

2) **The iterative process:** Given $\hat{\beta}^{(t)}$, $\hat{\sigma}_\varepsilon^{2(t)}$, $\hat{u}^{(t)} = (\hat{u}_1^{(t)T}, \dots, \hat{u}_m^{(t)T})^T$, and $\hat{D}^{(t)}$ from iteration t , compute estimates and predictions at iteration $t+1$ as follows:

a) $\hat{A}_j^{(t+1)} = Z_j^T Z_j + \hat{\sigma}_\varepsilon^{2(t)} [\hat{D}^{(t)}]^{-1}$, $j = 1, \dots, m$;

b) Solve for $\hat{\beta}^{(t+1)}$ and $\hat{u}^{(t+1)}$ in the linear system, with $\hat{A}^{(t+1)} = \text{diag}(\hat{A}_1^{(t+1)}, \dots, \hat{A}_m^{(t+1)})$:

$$\begin{pmatrix} X^T X & X^T Z \\ Z^T X & Z^T Z + \hat{A}^{(t+1)} \end{pmatrix} \begin{pmatrix} \hat{\beta}^{(t+1)} \\ \hat{u}^{(t+1)} \end{pmatrix} = \begin{pmatrix} X^T y \\ Z^T y \end{pmatrix}; \tag{29}$$

c) $\hat{\sigma}_\varepsilon^{2(t+1)} = y^T (y - X^{(t+1)} - Z\hat{u}^{(t+1)}) / (n - p)$;

d) $\hat{V}_j^{*(t+1)} = \hat{\sigma}_\varepsilon^{2(t+1)} [\hat{A}_j^{(t+1)}]^{-1}$, $j = 1, \dots, m$;

e) $\hat{D}^{(t+1)} = \frac{1}{m} \sum_{j=1}^m [\hat{V}_j^{*(t+1)} + \hat{u}_j^{(t+1)} \hat{u}_j^{(t+1)T}]$;

3) **Stopping criterion:** Assume convergence when

$\|\hat{D}^{(t+1)} - \hat{D}^{(t)}\|^{(M)} \leq \delta_1 \|\hat{D}^{(t+1)}\|^{(M)}$, $\|\hat{\beta}^{(t+1)} - \hat{\beta}^{(t)}\| \leq \delta_2 \|\hat{\beta}^{(t+1)}\|$, $|\hat{\sigma}_\varepsilon^{2(t+1)} - \hat{\sigma}_\varepsilon^{2(t)}| \leq \delta_3 \hat{\sigma}_\varepsilon^{2(t+1)}$ are all satisfied, where $\delta_1, \delta_2, \delta_3$ are relative tolerance levels set in $(0, 1)$, and $\|\cdot\|^{(M)}$ is a chosen matrix norm. Otherwise, $t \leftarrow t+1$ and repeat Step 2.

4) **Extracting estimates:** At convergence, take $\hat{\beta} = \hat{\beta}^{(t+1)}$, $\hat{\sigma}_\varepsilon^2 = \hat{\sigma}_\varepsilon^{2(t+1)}$, $\hat{D} = \hat{D}^{(t+1)}$ as estimates of β , σ_ε^2 , D . Also, take $\hat{u}_1^{(t+1)}, \dots, \hat{u}_m^{(t+1)}$ as predictions of U_1, \dots, U_m .

Remark 8. In that algorithm, whereas stabilization of the D and σ_ε^2 iterates is enough to ensure that the values of all the other estimates have stabilized as well, the stopping criterion of the algorithm is based on also monitoring the β iterates for more security.

Given the alternative system of EEs (28a)-(28c), our second algorithm for 2-level LMMs with u.h.o. errors goes as follows:

Algorithm 3S-A1-V2. Estimating β , σ_ε^2 , D and predicting U_1, \dots, U_m in a 2-level LMM with u.h.o. errors. Version 2

1) **Initialization:** The same as in Algorithm 3S-A1-V1;
 2) **The iterative process:** Given $\hat{\beta}^{(t)}$, $\hat{\sigma}_\varepsilon^{2(t)}$, $\hat{\mathbf{u}}_1^{(t)}, \dots, \hat{\mathbf{u}}_m^{(t)}$, hence $\hat{\mathbf{u}}^{(t)} = (\hat{\mathbf{u}}_1^{(t)\top}, \dots, \hat{\mathbf{u}}_m^{(t)\top})^\top$, and $\hat{\mathbf{D}}^{(t)}$ from iteration t , compute estimates and predictions at iteration $t+1$ as follows:

$$\begin{aligned} \text{a) } \hat{\mathbf{B}}_j^{(t+1)} &= \left(\mathbf{Z}_j^\top \mathbf{Z}_j + \hat{\sigma}_\varepsilon^{2(t)} \left[\hat{\mathbf{D}}^{(t)} \right]^{-1} \right)^{-1}, \quad j=1, \dots, m; \\ \text{b) } \hat{\beta}^{(t+1)} &= \left(\mathbf{X}^\top \mathbf{X} \right)^{-1} \mathbf{X}^\top \left(\mathbf{y} - \mathbf{Z} \hat{\mathbf{u}}^{(t)} \right); \\ \text{c) } \hat{\mathbf{u}}_j^{(t+1)} &= \hat{\mathbf{B}}_j^{(t+1)} \mathbf{Z}_j^\top \left(\mathbf{y}_j - \mathbf{X}_j \hat{\beta}^{(t+1)} \right), \quad j=1, \dots, m; \\ \text{d) } \hat{\sigma}_\varepsilon^{2(t+1)} &= \frac{1}{n-p} \left(\left\| \mathbf{y} - \mathbf{Z} \hat{\mathbf{u}}^{(t+1)} - \mathbf{X} \hat{\beta}^{(t+1)} \right\|^2 + \hat{\sigma}_\varepsilon^{2(t)} \left\| \hat{\mathbf{u}}^{(t+1)} \right\|_{\hat{\mathbf{G}}^{(t-1)}}^2 \right), \\ \hat{\mathbf{G}}^{(t)} &= \text{diag} \left(\hat{\mathbf{D}}^{(t)}, \dots, \hat{\mathbf{D}}^{(t)} \right); \\ \text{e) } \hat{\mathbf{V}}_j^{*(t+1)} &= \hat{\sigma}_\varepsilon^{2(t+1)} \hat{\mathbf{B}}_j^{(t+1)}, \quad j=1, \dots, m; \\ \text{f) } \hat{\mathbf{D}}^{(t+1)} &= \frac{1}{m} \sum_{j=1}^m \left[\hat{\mathbf{V}}_j^{*(t+1)} + \hat{\mathbf{u}}_j^{(t+1)} \hat{\mathbf{u}}_j^{(t+1)\top} \right]; \end{aligned}$$

3) **Stopping criterion and extracting estimates.** As in Algorithm 3S-A1-V1.

Remark 9. Important elements about the practical implementation of the above two designed algorithms are detailed further in Section S2 of the accompanying Supplementary material document. In particular, it is explained there how we monitor a possible rank deficiency of \mathbf{D} and even whether or not random effects are really there for a given data set, in the first place. For the latter aspect, we introduce and motivate

$$\rho = \left\| \mathbf{Z} \mathbf{G} \mathbf{Z}^\top \right\|^{(M)} / \sigma_\varepsilon^2, \quad (30)$$

a quantity called random effects ratio, through formula (S:2.6b) to try to assess the likelihood of having significant random effects or not in the data.

4.3. First Properties of the Estimates $\hat{\beta}$, $\hat{\sigma}_\varepsilon^2$, $\hat{\mathbf{D}}$ and $\hat{\mathbf{u}}$

The convergence of the above two algorithms has not been investigated yet. However, the fact that their convergence can only occur to estimates $\hat{\beta}$, $\hat{\sigma}_\varepsilon^2$, $\hat{\mathbf{u}}$, $\hat{\mathbf{D}}$ (of β , σ_ε^2 , \mathbf{D} , $\bar{\mathbf{u}}$) satisfying the two equivalent systems of EEs of Section 0.0.5 allows to identify some first properties of those estimates. Indeed, assume that any of the above two algorithms has converged to $\hat{\beta}$, $\hat{\sigma}_\varepsilon^2$, $\hat{\mathbf{D}}$, $\hat{\mathbf{u}}$, with $\hat{\mathbf{D}}$ SPD, $\hat{\sigma}_\varepsilon^2 > 0$ and Assumption $\mathcal{A}2_g$ true. We highlight two important properties of these estimates.

4.3.1. Relationship with the HMMEs

Firstly, Theorems 1 & 4 and Corollary 1 imply the following about $\hat{\beta}$, $\hat{\sigma}_\varepsilon^2$, $\hat{\mathbf{D}}$, $\hat{\mathbf{u}}$:

Theorem 6. If one takes $\mathbf{G} = \hat{\mathbf{G}}$ and $\sigma_\varepsilon^2 = \hat{\sigma}_\varepsilon^2$ in the HMMEs (10), then $\tilde{\beta} = \hat{\beta}$ and $\tilde{\mathbf{U}} = \hat{\mathbf{u}}$ are the solutions to (10). Therefore, the estimates $\hat{\beta}$, $\hat{\sigma}_\varepsilon^2$, $\hat{\mathbf{D}}$ and prediction $\hat{\mathbf{u}}$ satisfy:

$$\hat{\beta} = \text{EBLUE} \left(\beta \mid \mathbf{y}, \hat{\mathbf{G}}, \hat{\sigma}_\varepsilon^2 \right) = \left(\mathbf{X}^\top \hat{\mathbf{V}}^{-1} \mathbf{X} \right)^{-1} \mathbf{X}^\top \hat{\mathbf{V}}^{-1} \mathbf{y}, \quad (31)$$

$$\hat{u} = \text{EBLUP}(U | y, \hat{G}, \hat{\sigma}_\varepsilon^2) = \hat{GZ}^T \hat{V}^{-1} (y - X \hat{\beta}), \tag{32}$$

where $\hat{V} = Z\hat{GZ}^T + \hat{\sigma}_\varepsilon^2 I_n$. That is $\hat{\beta}$ and \hat{u} are, respectively, an empirical BLUE of β and an empirical BLUP of U when G is estimated by \hat{G} and σ_ε^2 by $\hat{\sigma}_\varepsilon^2$.

Remark 10. The main difference between Algorithms 3S-A1-V1 and 3S-A1-V2 is that the former enforces satisfaction of the HMMEs for the computed estimates at each iteration $t \geq 1$, while the latter only guarantees that for the final estimates if convergence is achieved.

4.3.2. Relationship with the 2-Level Gaussian LMM Likelihood Equations under i.i.d. Errors

The EEs (27b) and (28b) indicate that our 2 estimation procedures above are extensions to 2-level LMMs with u.h.o. errors of the well known 2-step method to fit an LM with such errors whereby β is first estimated by OLS and, then, that estimator is used to get an unbiased estimator for σ_ε^2 , and this without using any parametric distributional assumption. Now, it is trivial that for the LM with i.i.d. Gaussian errors, that 2-step estimation method is asymptotically equivalent to the ML estimation of β and σ_ε^2 . We are going to show that the same property holds for our 2 estimation procedures when the 2-level LMM is, indeed, a Gaussian one with i.i.d. errors, at least in terms of the solutions $\hat{\beta}$, $\hat{\sigma}_\varepsilon^2$, \hat{D} satisfying the *likelihood equations*. So, we assume a 2-level LMM satisfying Assumptions $\mathcal{A}1 - \mathcal{A}4$ and, for $j=1, \dots, m$:

Assumption $\mathcal{A}5$. $U_j \stackrel{\mathcal{L}}{\sim} \mathcal{N}_r(\mathbf{0}, \mathbf{D})$ and $\varepsilon_j \stackrel{\mathcal{L}}{\sim} \mathcal{N}_{n_j}(\mathbf{0}, \sigma_\varepsilon^2 I_{n_j})$.

Estimating the parameters β , σ_ε^2 , \mathbf{D} of that 2-level Gaussian LMM through ML usually starts by trying to solve the likelihood equations:

$$\frac{\partial}{\partial \beta} \ell(\beta, \sigma_\varepsilon^2, \mathbf{D} | \mathbf{y}) = \mathbf{0}, \tag{33a}$$

$$\frac{\partial}{\partial \mathbf{D}} \ell(\beta, \sigma_\varepsilon^2, \mathbf{D} | \mathbf{y}) = \mathbf{0}, \tag{33b}$$

$$\frac{\partial}{\partial \sigma_\varepsilon^2} \ell(\beta, \sigma_\varepsilon^2, \mathbf{D} | \mathbf{y}) = 0, \tag{33c}$$

where $\ell(\beta, \sigma_\varepsilon^2, \mathbf{D} | \mathbf{y})$ is the log-likelihood of the model with observed response $\mathbf{Y} = \mathbf{y}$.

First, it is well known that, given our assumptions:

$$(33a) \Leftrightarrow \beta = (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1} \mathbf{y}. \tag{34}$$

Thus, thanks to (31), the triplet $(\beta, \sigma_\varepsilon^2, \mathbf{D}) = (\hat{\beta}, \hat{\sigma}_\varepsilon^2, \hat{D})$ satisfies the first likelihood Equation (33a). For the other two likelihood equations, generalizing computations in Searle *et al.* ([14], pages 278-279) for the ANOVA model (4), one has (Proof in the **Appendix**):

Theorem 7. Let $\tilde{\beta}$ and \tilde{U} satisfy the HMMEs (10) given σ_ε^2 , \mathbf{D} for the 2-level LMM (2), with Assumptions $\mathcal{A}1 - \mathcal{A}4$, $\mathcal{A}2_g$ and $\mathcal{A}5$.

1) If (33a) holds, then:

$$(33b) \Leftrightarrow \mathbf{D} = \frac{1}{m} \sum_{j=1}^m (\mathbf{V}_j^* + \tilde{U}_j \tilde{U}_j^T).$$

2) If (33a) and (33b) both hold, then:

$$(33c) \Leftrightarrow \sigma_\varepsilon^2 = \mathbf{y}^T (\mathbf{y} - \mathbf{X} \tilde{\beta} - \mathbf{Z} \tilde{U}) / n.$$

Given the EE (27c), the first equivalence in that theorem shows that the second likelihood Equation (33b) is satisfied by the quadruplet

$(\beta, \sigma_\varepsilon^2, \mathbf{D}, U) = (\hat{\beta}, \hat{\sigma}_\varepsilon^2, \hat{\mathbf{D}}, \hat{U})$. But, given our EE (27b), we see that the triplet

$(\beta, \sigma_\varepsilon^2, U) = (\hat{\beta}, \hat{\sigma}_\varepsilon^2, \hat{U})$ falls short to satisfying (33c) only through having a denominator $n - p$ instead of n , exactly like in the LM case. So, when the 2-level LMM is a Gaussian one with i.i.d. errors, our 2 estimation procedures above are asymptotically equivalent to the ML method. However, as for the LM case, we stick here with our denominator $n - p$, inspired by the unbiasedness property (16a) in Theorem 5.

Remark 11. *An important by-product of Theorem 7 is that it implies that by replacing $n - p$ by n in the denominator of the formula updating the residual variance estimate during the iterations in Algorithms 3S-A1-V1 and 3S-A1-V2, we get two new algorithms for computing the Gaussian ML estimates in a 2-level LMM. They may be coded 3S-ML-A1-V1 and 3S-ML-A1-V2 and can be viewed as alternatives to existing algorithms based on Newton-Raphson, Fisher scoring or EM (Expectation-Maximization).*

4.4. 3S LMM Fitting: What about the Accuracy of the Estimates?

Obviously, once a 3S algorithm has converged to the targeted estimates of parameters for fitting an LMM to a given data set, the next question is, as usual in statistical inference: *what about the accuracy of those estimates?* In a strictly parametric modeling with Gaussian ML or REML estimation, the answer is customarily built through estimating the inverse of Fisher's Information Matrix for the estimated parameters. But we cannot go that route here because we precisely aimed here at a nonparametric modeling with the LMM for the provided data set. Instead, to assess the accuracy of the estimates computed by a 3S algorithm, we use the bootstrap. But bootstrapping mixed models necessitates significantly more care than is usually the case when using routine i.i.d. data. Methods and up-to-date discussions about this can be found for instance in Carpenter *et al.* [45], Van der Leeden *et al.* [46], Chambers and Chandra [47], Thai *et al.* [48], Modugno and Giannerini [49].

5. 3S Fitting of LMMs: Beyond 2-Level LMMs

As stated from the outset, the first goal in this work was to design a methodology for parameters estimation in a 2-level LMM where the only assumption added to the basic ones (\mathcal{A}_1 to \mathcal{A}_4) would be to have u.h.o. errors (Assumption \mathcal{A}_g). Nonetheless, analyzing the 3-step sequence upon which each of our devised 3S iterative algorithms was constructed in Section 4, one is struck by the

fact that the 2-level element only intervenes in Step 3, *i.e.* when deriving an EE for \mathbf{D} . In the construction of both algorithms, what is done in Steps 1-2 is based on results derived in Section 3.2.2 and classical formulas of the BLP in an LMM and its covariance matrix, thus only uses Assumptions \mathcal{A}_{1_g} , \mathcal{A}_{2_g} and \mathcal{A}_{4_g} of an arbitrary LMM (1). It comes that for any LMM satisfying the latter 3 assumptions, if one can reliably devise an EE for \mathbf{G} , then one would immediately get respective adaptations of Algorithms 3S-A1-V1 and 3S-A1-V2 for fitting it. In that order, one can start from the fact that (analogous to (18) for the 2-level case)

$$\tilde{\mathbf{G}}_1 = \mathbf{V}^* + \bar{\mathbf{U}}\bar{\mathbf{U}}^T \tag{35}$$

is an unbiased estimator of \mathbf{G} , where $\mathbf{V}^* = \mathbf{G} - \mathbb{M}\text{cov}(\bar{\mathbf{U}})$.

The most obvious case is when \mathbf{G} is a diagonal matrix, but with some diagonal elements being equal by design. We examine two classical such situations hereafter: the 2-level LMM with \mathbf{D} diagonal and the ANOVA LMM (4).

5.1. 3S Fitting of 2-Level LMMs with u.ho. Errors and \mathbf{D} Diagonal

In a 2-level LMM (2), when the dimension r of the vectors U_1, \dots, U_m is big while the dimension n of the response vector is small to moderate, the number $p + mr + r(r+1)/2$ of scalar parameters to estimate in β , \mathbf{D} and the random effects predictors might become too big for any LMM fitting method. In that case, imposing a diagonal structure for \mathbf{D} may be a convenient way to achieve a reasonable estimation process for β , \mathbf{D} and σ_ε^2 from the observed response vector $\mathbf{Y} = \mathbf{y} \in \mathbb{R}^n$. This motivates assuming then that $\mathbf{D} = \text{diag}(\sigma_1^2, \dots, \sigma_r^2)$, with $\sigma_k^2 > 0, \forall k = 1, \dots, r$.

With that, the unbiased estimator $\tilde{\mathbf{D}}_1$ of \mathbf{D} given by (18) implies that for each $k = 1, \dots, r$, an unbiased preliminary estimator of the variance component σ_k^2 is

$$\tilde{\sigma}_k^2 = \frac{1}{m} \sum_{j=1}^m \left[(\mathbf{V}_j^*)_{kk} + (\bar{U}_j)_k^2 \right], \tag{36}$$

where $(\mathbf{V}_j^*)_{kk}$ and $(\bar{U}_j)_k$ are the k^{th} diagonal element of matrix \mathbf{V}_j^* and k^{th} component of vector \bar{U}_j . Thus the EE (21) for \mathbf{D} is simplified here to:

$$\hat{\mathbf{D}} = \text{diag}(\hat{\sigma}_1^2, \dots, \hat{\sigma}_r^2), \quad \hat{\sigma}_k^2 = \frac{1}{m} \sum_{j=1}^m \left[(\hat{\mathbf{V}}_j^*)_{kk} + (\hat{\mathbf{u}}_j)_k^2 \right]. \tag{37}$$

Then one can easily adapt Algorithms 3S-A1-V1 and 3S-A1-V2 for the 2-level LMM with u.ho. errors and diagonal \mathbf{D} , to get, say, Algorithms 3S-A1-V1-diag and 3S-A1-V2-diag. The latter are presented in Section S4.1 of the Supplementary material document.

5.2. 3S Fitting of an ANOVA LMM

In the ANOVA LMM (4), we can partition $\bar{\mathbf{U}} = \text{BLP}(U | \mathbf{Y})$ as in (S:1.6a), but where each \bar{U}_j is rather a q_j -vector ($j = 1, \dots, m$). Now, let: $s_0 = 0$, and

$\forall j = 1, \dots, m$, $s_j = q_1 + \dots + q_j = s_{j-1} + q_j$. Then each vector \bar{U}_j is made up of the $(s_{j-1} + 1)^{\text{th}}$ to $(s_j)^{\text{th}}$ components of \bar{U} .

With the structure of \mathbf{G} in this case, the estimator $\tilde{\mathbf{G}}_1$ in (35) implies that an unbiased preliminary estimator of each variance component σ_j^2 is

$$\tilde{\sigma}_j^2 = \frac{1}{q_j} \left[\sum_{k=s_{j-1}+1}^{s_j} (\mathbf{V}^*)_{kk} + \|\bar{U}_j\|^2 \right] \quad (38)$$

for $j = 1, \dots, m$. Hence, we get an EE for \mathbf{G} here as:

$$\hat{\mathbf{G}} = \text{diag}(\hat{\sigma}_1^2, \dots, \hat{\sigma}_m^2), \quad \hat{\sigma}_j^2 = \frac{1}{q_j} \left[\sum_{k=s_{j-1}+1}^{s_j} (\hat{\mathbf{V}}^*)_{kk} + \|\hat{\mathbf{u}}_j\|^2 \right]. \quad (39)$$

From that, Algorithms 3S-A1-V1 and 3S-A1-V2 can be adapted to fit an ANOVA LMM to get, say, Algorithms 3S-A1-V1-ANOVA and 3S-A1-V2-ANOVA. They are presented in Section S4.2 of the Supplementary material. The first one is nearly the same as one of those derived by Henderson for Gaussian ANOVA LMMs [14] [16]. The only difference is again the $n - p$ in the denominator in the r.h.s. of the EE for σ_ε^2 in (27b).

6. Numerical Examples

In this section, we carry out some numerical experiments on simulated data sets (Section 6.1) and two real world data sets (Section 6.2). For both situations, we compared Algorithm 3S-A1-V1 vs. Gaussian ML. For the latter, and in light of Theorem 7, we could have used Algorithm 3S-ML-A1-V1 in Remark 11. But that would have looked like a self comparison. So we instead chose to use the implementation of Gaussian ML by the `lmer` function in the reference `lme4` package [50] of the R software. Hereafter, we will denote the latter by `lmer-ML`. The reason why Algorithm 3S-A1-V2 is not included in the study is that, as expected, it tends to give the same results as Algorithm 3S-A1-V1.

6.1. A Simulation Study

Here, to investigate the performance of Algorithm 3S-A1-V1 vs. `lmer-ML`, we fitted both to 500 simulated data sets, each of size $n = 200$, under various distributional scenarios. First, we assumed we have a population Ω comprising $m = 10$ clusters in equal proportions of $1/10$ and that in each cluster j , both the vector X of fixed effects covariates and Z , that of random effects covariates, follow Gaussian distributions with given parameters, respectively in \mathbb{R}^4 and \mathbb{R}^3 . To get a unit in a simulated data set, we first sample its cluster among the 10 in Ω , then its covariates by sampling from the Gaussian distributions of X and Z in that cluster. The distributions parameters used for that and the simulated distributions described below are given in Section S5 of the Supplementary material document, alongside what we used as vector of fixed effects parameters $\beta \in \mathbb{R}^4$.

To get one data set, after simulating the 200 units in it with their fixed and random effects covariates as just described,

- we first simulate a vector of random effects $U_j \in \mathbb{R}^3$ for each cluster

- $j \in \{1, 2, \dots, 10\}$ according to a chosen distribution;
- then, for each unit i sampled in cluster j and included in the data set, with simulated fixed effects covariates $X_{ij} \in \mathbb{R}^4$ and random effects covariates $Z_{ij} \in \mathbb{R}^3$, and a residual error $\varepsilon_{ij} \in \mathbb{R}$ according to a chosen distribution. Then its simulated exact response is calculated as $Y_{ij} = X_{ij}^T \beta + Z_{ij}^T U_j + \varepsilon_{ij}$.

For the distribution of the residual errors, we examined three options: a Gaussian, a mixture of 4 Gaussians with given mixing proportions, and a discrete distribution with 4 mass points. On the other hand, the covariance matrix D of each U_j is 3×3 . For the distribution of the U_j 's in \mathbb{R}^3 , we examined five options: Gaussian with D SPD, Gaussian with $\text{rank}(D) = 1$, Gaussian with $\text{rank}(D) = 2$, mixture of 3 Gaussians with given mixing proportions, and a discrete distribution with 3 mass points. Thus, we have $3 \times 5 = 15$ possible scenarios. Hence, for each of the algorithms 1 and lmer-ML, in each of the 15 scenarios, we simulated 500 200-sized data sets and fitted the 2-level LMM (expressed in the R software notations):

$$Y \sim \text{fe}(-1 + X1 + X2 + X3 + X4) + \text{re}(-1 + Z1 + Z2 + Z3) + \text{Gr}(\text{cluster}) \quad (40)$$

where “-1” signifies no intercept included, be it on the fixed or the random effects.

The fit of the LMM (40) to each simulated data set yields an estimate $\hat{\theta} \in \{\hat{\beta}, \hat{D}, \hat{\sigma}^2\}$ of each parameter $\theta \in \{\beta, D, \sigma^2\}$, and a prediction $\hat{U} = (\hat{U}_j^T)_{j=1, \dots, 10}^T$ of the random effects vector $U = (U_j^T)_{j=1, \dots, 10}^T \in \mathbb{R}^{30}$.

Furthermore, to estimate the Mean Squared Prediction Error (MSPE) of the unit response Y by its prediction \hat{Y} from a fitted model, for each simulated data, we independently simulated an additional data set of size 100 with simulated response vector $Y \in \mathbb{R}^{100}$ and computed \hat{Y} , the response vector predicted by the fitted LMM. For each $\theta \in \{\beta, D, \sigma^2\}$, we then used its simulated replicates $\hat{\theta}_1, \dots, \hat{\theta}_{500}$ to estimate

$$b_R(\hat{\theta}) = \|\mathbb{E}(\hat{\theta}) - \theta\| / \|\theta\| \text{ and } \text{RMSE}(\hat{\theta}) = \mathbb{E} \|\hat{\theta} - \theta\|^2 / \|\theta\|^2, \quad (41)$$

respectively the Euclidean norm of the Relative Bias and the Relative Mean Squared Error of $\hat{\theta}$. Each expectation was estimated by the corresponding Monte Carlo estimate. In that calculation, for the symmetric 3×3 matrices $\theta = D$ and $\hat{\theta} = \hat{D}$, we identified, each, to the vector of elements in its upper triangular part (comprising the diagonal). Likewise, for each predictor $\hat{P} \in \{\hat{U}, \hat{Y}\}$ of $P \in \{U, Y\}$, we used its simulated replicates to estimate the norm of its bias and its MSPE.

$$b(\hat{P}) = \|\mathbb{E}(\hat{P}) - P\| \text{ and } \text{MSPE}(\hat{P}) = \mathbb{E} \|\hat{P} - P\|^2. \quad (42)$$

For both algorithms, the results about the estimation of β , D , σ_ε^2 and the prediction of random effects and responses are reported in **Table 1**. Before we specifically comment on them, note, however, that when running an iterative algorithm on simulated data sets, the algorithm might fail to converge on some

Table 1. Estimates of the bias and MSE of estimates of parameters, predicted random effects and responses in the LMM (40) applied to simulated data sets with different distributions for the random effects and residual errors.

Distribution of random effects and errors	$b_R(\hat{\beta})$	$\sqrt{\text{RMSE}(\hat{\beta})}$	$b_R(\hat{\sigma}_\varepsilon^2)$	$\sqrt{\text{RMSE}(\hat{\sigma}_\varepsilon^2)}$	$b_R(\hat{\mathbf{D}})$	$\sqrt{\text{RMSE}(\hat{\mathbf{D}})}$	$b(\hat{U})$	$\sqrt{\text{MSPE}(\hat{U})}$	$b(\hat{Y})$	$\sqrt{\text{MSPE}(\hat{Y})}$
Gaussian errors										
Gaussian RE (\mathbf{D} SPD)	0.00355 (0.000861)	0.0689 (0.0678)	0.00216 (0.0188)	0.112 (0.108)	0.532 (0.521)	0.958 (0.951)	0.0165 (0.00782)	0.662 (0.673)	$3.80e-05$ ($9.46e-05$)	$5.96e-03$ ($5.87e-03$)
Gaussian RE ($\mathbf{R}(\mathbf{D}) = 1$)	0.00327 (0.00351)	0.0535 (0.0540)	0.0172 (0.0410)	0.103 (0.105)	0.878 (0.918)	1.02 (2.39)	0.00693 (0.00202)	0.358 (0.363)	$4.55e-06$ ($2.18e-05$)	$5.17e-03$ ($5.28e-03$)
Gaussian RE ($\mathbf{R}(\mathbf{D}) = 2$)	0.00449 (0.00177)	0.0625 (0.0624)	0.0164 (0.0401)	0.108 (0.115)	0.74 (0.765)	1.26 (1.25)	0.00350 (0.00291)	0.533 (0.551)	$4.20e-05$ ($7.57e-05$)	$5.76e-03$ ($5.66e-03$)
Gaussian mixture RE	0.0022 (0.00152)	0.0676 (0.0665)	0.0055 (0.0155)	0.109 (0.104)	0.425 (0.521)	1.15 (0.737)	0.0116 (0.00709)	0.684 (0.674)	$1.45e-04$ ($6.95e-05$)	$4.43e-03$ ($4.43e-03$)
Discrete RE	0.00224 (0.00249)	0.0606 (0.0612)	0.0219 (0.0380)	0.105 (0.111)	0.818 (0.718)	1.84 (1.21)	0.00411 (0.00553)	0.513 (0.516)	$2.91e-05$ ($6.56e-05$)	$3.57e-03$ ($3.63e-03$)
Gaussian mixture errors										
Gaussian RE (\mathbf{D} SPD)	0.00109 (0.00254)	0.0641 (0.0644)	0.364 (0.416)	0.406 (0.446)	0.724 (0.567)	1.44 (1.12)	0.00819 (0.00515)	0.598 (0.598)	$9.23e-05$ ($8.87e-05$)	$5.84e-03$ ($5.88e-03$)
Gaussian RE ($\mathbf{R}(\mathbf{D}) = 1$)	0.00165 (0.00421)	0.0622 (0.0621)	0.404 (0.457)	0.441 (0.483)	0.935 (2.87)	2.20 (7.43)	0.00377 (0.00484)	0.438 (0.461)	$2.21e-05$ ($5.48e-05$)	$5.42e-03$ ($5.16e-03$)
Gaussian RE ($\mathbf{R}(\mathbf{D}) = 2$)	0.00368 (0.00300)	0.0628 (0.0622)	0.402 (0.417)	0.441 (0.452)	0.670 (0.816)	0.991 (1.01)	0.00212 (0.00764)	0.519 (0.532)	$2.80e-05$ ($1.48e-04$)	$5.72e-03$ ($5.65e-03$)
Gaussian mixture RE	0.00375 (0.00287)	0.0626 (0.0646)	0.395 (0.404)	0.433 (0.438)	0.510 (0.722)	0.898 (0.786)	0.00826 (0.00707)	0.609 (0.596)	$1.14e-06$ ($2.58e-05$)	$4.45e-03$ ($4.43e-03$)
Discrete RE	0.00331 (0.00310)	0.0641 (0.0632)	0.384 (0.426)	0.424 (0.457)	0.728 (0.638)	1.18 (1.25)	0.00756 (0.00637)	0.512 (0.509)	$1.87e-05$ ($2.75e-06$)	$3.69e-03$ ($3.75e-03$)
Discrete errors										
Gaussian RE (\mathbf{D} SPD)	0.00134 (0.00195)	0.0577 (0.0552)	0.774 (0.784)	0.778 (0.788)	0.643 (0.589)	1.57 (1.37)	0.00341 (0.00724)	0.507 (0.504)	$6.21e-05$ ($2.56e-05$)	$5.99e-03$ ($5.80e-03$)
Gaussian RE ($\mathbf{R}(\mathbf{D}) = 1$)	0.00132 (0.00301)	0.0597 (0.0608)	0.778 (0.792)	0.782 (0.796)	0.723 (0.737)	1.43 (1.73)	0.00364 (0.00410)	0.472 (0.479)	$1.23e-05$ ($7.19e-05$)	$5.20e-03$ ($5.19e-03$)
Gaussian RE ($\mathbf{R}(\mathbf{D}) = 2$)	0.00221 (0.00266)	0.0586 (0.0595)	0.773 (0.783)	0.778 (0.787)	1.10 (1.51)	2.55 (3.27)	0.00726 (0.00719)	0.486 (0.490)	$5.39e-07$ ($3.26e-05$)	$5.73e-03$ ($5.76e-03$)
Gaussian mixture RE	0.00174 (0.00266)	0.0570 (0.0560)	0.778 (0.777)	0.783 (0.781)	0.384 (0.614)	0.955 (1.39)	0.00119 (0.00358)	0.500 (0.510)	$7.66e-05$ ($1.20e-06$)	$4.53e-03$ ($4.53e-03$)
Discrete RE	0.00188 (0.00168)	0.0599 (0.0598)	0.772 (0.786)	0.777 (0.790)	0.848 (0.633)	1.40 (1.35)	0.00786 (0.00666)	0.485 (0.489)	$2.97e-05$ ($1.60e-05$)	$3.69e-03$ ($3.62e-03$)

\mathbf{D} is the cluster random effects (RE) 3×3 covariance matrix; $\mathbf{R}(\mathbf{D}) = \text{rank}(\mathbf{D})$ is the rank of \mathbf{D} . In each cell of two numbers in the table, the one on top is the result from the Algorithm 3S-A1-V1, while the one below it, in parentheses, is the result from Gaussian ML.

of them. In our case, the strategy was to discard such data sets, but, in each scenario, run the simulation until 500 data sets have been successfully fitted by the considered algorithm. Nonetheless, while Algorithm 3S-A1-V1 never failed during this set of simulations, the situation appeared particularly tricky when running Algorithm lmer-ML. Indeed, with its default settings, its inner workings make that it signaled a failure of convergence on a massive number of simulated data sets, so that we decided actually not to reject them all, rejecting only the cases where the algorithm signaled a difficulty to converge due to a suspected singular fit, *i.e.* true model parameters likely on or close to the boundary of their space of values. Even with that, the only scenarios where we did not witness any failure of lmer-ML was under the Gaussian mixture distribution for the cluster random effects in \mathbb{R}^3 , whereas one would have only expected failure when the covariance matrix \mathbf{D} is rank deficient.

As might be expected, the results of both algorithms are comparable. Nevertheless, while there is no aspect in both tables in which the lmer-ML algorithm uniformly or almost uniformly dominates 3S-A1-V1, a dominance of the latter over the former can be observed, in almost all scenarios, on a lower bias of the estimate of σ_ε^2 , and not far from so on a lower MSE of that same estimate. Such an almost uniform dominance of the 3S-A1-V1 algorithm can also be seen in the

prediction of the response in terms of bias. As for comparing the two algorithms in various individual scenarios, the main striking observation is that 1 almost uniformly dominates lmer-ML when the cluster random effects covariance matrix \mathbf{D} is rank deficient, more decisively so the smaller that rank.

6.2. Two Real World Data Sets

In this section, we apply Algorithm 3S-A1-V1 on two classical data sets, cake [51] and Blackmore (from the car R package) data, and compare it to the results from Gaussian ML. The results of estimating β , σ_ε^2 , \mathbf{D} and predicting U_1, \dots, U_m by Algorithm 3S-A1-V1 will be presented. Each of these estimates is given with variability indicators such as estimates of Bias, Mean Squared Error (MSE) and t -statistics, p -values for significance and 95% two-sided confidence intervals (CI), all computed using the nonparametric residuals bootstrap approach outlined in Carpenter *et al.* [45]. In each LMM fit to a data set, we will also provide an estimate of the random effects ratio ρ in the data, given by (30).

6.3. Application to the Cake Data

The cake data set consists of observations of the *breakage angle* of chocolate cakes made with 3 different recipes and baked at 6 different temperatures, with 15 replicates for each combination of a recipe and a temperature. Hence, there are $3 \times 6 \times 15 = 270$ sample units (the baked cakes) for each of which there is a recorded value of the variables: angle (numeric), recipe (factor with levels A, B, C), temperature (ordered factor with levels $175 < 185 < 195 < 205 < 215 < 225$), replicate (factor with levels 1 to 15), temp (with the same values as temperature, but viewed as numeric). In the R software notations [39], we fitted to these data the two LMMs:

$$\text{angle} \sim \text{fe}(-1 + \text{temp}) + \text{re}(1) + \text{Gr}(\text{recipe} : \text{replicate}) \quad (43a)$$

$$\text{angle} \sim \text{fe}(-1 + \text{temperature}) + \text{re}(1) + \text{Gr}(\text{recipe} : \text{replicate}) \quad (43b)$$

In the LMM (43a), the response is angle, the fixed effects variable is temp, we have a random intercept as only random effects variable while the clustering variable is obtained by crossing the categories of the two factors recipe and replicate, thus yielding $3 \times 15 = 45$ clusters. The difference between the LMMs (43a) and (43b) is that the latter uses instead, as lone fixed effects variable, temperature viewed as a 6-level factor. The reason the intercept has been excluded in the fixed effects part of both models is that a preliminary fitting of the LMM (43a) with an intercept among the fixed effects revealed it to be insignificant through the bootstrap procedure. The interest in fitting the two models is that it allows to assess whether the true values of the temperature really matter in the quality of the fit.

The results for models (43a) and (43b) are respectively presented in **Table 2** and **Table 3**. For a given parameter estimate, the t -statistic is the estimate value divided by the square root of its MSE and the p -value is calculated assuming that t -statistic has, approximately, a standard Gaussian distribution when the true

Table 2. Results of algorithm 3S-A1-V1 fitting the LMM (43a) to the cake data.

parameters	est.	[2.5% – 97.5%]	bias	std.dev	$\sqrt{\text{MSE}}$	t-stat	p-value
β_{temp}	0.1604	[0.1519, 0.1705]	0.0002258	0.004798	0.004803	33.400	$1.545e - 244$
$\sigma_{\text{intercept}}^2$	39.28	[20.77, 60.00]	-0.1352	9.850	9.851	3.988	$6.675e - 05$
σ_{ε}^2	20.71	[17.01, 24.64]	-0.1051	1.946	1.949	10.62	$2.299e - 26$
ρ	1.897	[1.003, 2.978]	0.02226	0.5226	0.5231	3.626	$2.876e - 04$

β_{temp} : coeff. of fixed effects temp; $\sigma_{\text{intercept}}^2$: random intercept variance; σ_{ε}^2 : residual variance; ρ : random effects ratio; *est.*: estimate; [2.5% - 97.5%], bias, std.dev, $\sqrt{\text{MSE}}$, t-stat, p-value: respectively estimated two-sided 95% percentile confidence interval, bias, standard deviation, square root of the MSE, t-statistic and p-value of the t-test for a zero value, all computed by a bootstrap simulation of 1000 replicates of the model fit using the nonparametric residuals bootstrap approach described in Carpenter *et al.* [45]. These variability estimates are computed from the 1000 bootstrap replicates using the standard formulas introduced by Efron (see [52]).

Table 3. Results of Algorithm 3S-A1-V1 fitting the LMM (43b) to the cake data.

parameters	est.	[2.5% – 97.5%]	bias	std.dev	$\sqrt{\text{MSE}}$	t-stat	p-value
$\beta_{\text{temperature.175}}$	27.980	[25.89, 30.28]	0.01359	1.141	1.141	24.52	$8.184e - 133$
$\beta_{\text{temperature.185}}$	29.96	[27.82, 32.37]	-0.001430	1.176	1.176	25.47	$3.735e - 143$
$\beta_{\text{temperature.195}}$	31.42	[29.30, 33.92]	0.03457	1.191	1.192	26.36	$3.469e - 153$
$\beta_{\text{temperature.205}}$	32.18	[29.90, 34.51]	0.009740	1.194	1.194	26.94	$7.656e - 160$
$\beta_{\text{temperature.215}}$	35.84	[33.76, 38.43]	0.04854	1.189	1.190	30.13	$1.970e - 199$
$\beta_{\text{temperature.225}}$	35.36	[33.16, 37.83]	0.0009604	1.161	1.161	30.45	$1.121e - 203$
$\sigma_{\text{intercept}}^2$	39.30	[20.720, 60.00]	-0.2556	9.807	9.811	4.006	$6.168e - 05$
σ_{ε}^2	20.57	[16.66, 24.66]	-0.1119	2.042	2.045	10.06	$8.470e - 24$
ρ	1.911	[1.006, 3.040]	0.01894	0.5280	0.5284	3.616	$2.987e - 04$

Meaning of parameters as in **Table 2**.

value of the parameter is 0. The latter is rarely a bad approximation in this context. Overall, the two model fits are quite similar in terms of random effects variance, residual variance and MSPE estimate. Moreover, the relative difference in terms of fitted response values between the two fits ranges between -0.047 and 0.030 . Finally, the p-value of the random effects ratio ρ is roughly 3×10^{-4} , strongly suggesting the presence of random effects on the intercept w.r.t. the recipe:replicate clustering considered for the baked cakes.

For comparison, we also fitted each of the two LMMs (43a) and (43b) by the Gaussian ML method. The results are given in **Table 4** and, as expected, the parameters estimates are almost identical to those in the two previous tables.

6.4. Application to the Blackmore Longitudinal Data

To illustrate our methods on longitudinal LMMs, we consider the Blackmore data described in detail in Davis *et al.* [53]. It is a longitudinal, retrospective, self-reported data from a case-control study on the physical exercise histories of 231 teenage girls, with 138 who are eating disorder (*anorexia nervosa*) patients recruited from a 4 years inpatient Eating Disorder Program at the Toronto hospital for sick children. The other 93 comparable “control” subjects with no history of a psychiatric disorder (determined by asking relevant clinical questions) were recruited from informational letters through school boards to parents, inviting the teenage daughter (the letter bearer) to take part in the study. Retrospective recall of leisure-time sport and exercise activities at target ages 8, 10, 12, 14 years (if applicable in relation to the girl’s current age), and the 12 months

Table 4. LMM Gaussian maximum likelihood fitting to the cake data.

				LMM (43b)			
				<i>parameters</i>	<i>est.</i>	<i>std.dev</i>	<i>t-stat</i>
LMM (43a)							
<i>parameters</i>	<i>est.</i>	<i>std.dev</i>	<i>t-stat</i>				
β_{temp}	0.1604	0.004666	34.38	$\beta_{temperature.175}$	27.98	1.149	24.35
$\sigma_{intercept}^2$	39.30			$\beta_{temperature.185}$	29.96	1.149	26.07
σ_{ϵ}^2	20.62			$\beta_{temperature.195}$	31.42	1.149	27.35
				$\beta_{temperature.205}$	32.18	1.149	28.00
				$\beta_{temperature.215}$	35.84	1.149	31.19
				$\beta_{temperature.225}$	35.36	1.149	30.77
				$\sigma_{intercept}^2$	39.40		
				σ_{ϵ}^2	20.02		

prior to the study, were obtained during a structured interview. Since the girls were recruited at different ages, the number of observations and the age at the last observation vary, respectively from 2 to 5 observations and from 11.58 to 17.92 years. The Blackmore data are given in the long format for a longitudinal data set, with the following variables: subject (an identification code for each girl); age (the girl's age in years at the time of observation); exercise (the variable of interest, the amount of physical activity in which the girl engaged, expressed as estimated hours per week) and group (a factor indicating whether the girl is a "patient" or a "control"). For modeling exercise in terms of age and group across the sample of girls, Davis *et al.* [53] first \log_2 -transformed exercise to make its distribution for both groups more symmetric and linearised its relationship with age. Because there are some 0 values of exercise, 5 minutes (5/60 of an hour) were added to each value of exercise prior to taking logs.

For these data, we fit an LMM using as fixed effects: an intercept, age minus 8 (denoted age8), group and the interaction of the two latter. As for random effects in this study, they are attached to variations from girl to girl, so the clustering variable is subject. Now, a follow-up plot of \log_2 exercise by age for 20 randomly selected patients and 20 randomly selected control girls showed that the \log_2 exercise at the start of follow-up varies considerably, suggesting a random intercept in the LMM. Also, the evolution of \log_2 exercise with age differs from girl to girl, informing on the inclusion of a random slope for age8. The model to fit is therefore:

$$\log_2 \text{ exercise} \sim \text{fe}(1 + \text{age8} + \text{group} + \text{age8} * \text{group}) + \text{re}(1 + \text{age8}) + \text{Gr}(\text{subject}) \quad (44)$$

when fitting that model, and as traditional in this type of study, in the numerical coding of the binary variable group, our reference category is "control". Moreover, the covariance matrix \mathbf{D} is 2×2 , with diagonal elements $\sigma_{intercept}^2$ and σ_{age8}^2 , and both off-diagonal ones equal to $\text{cov}_{intercept, age8}$.

The results for the model using Algorithm 3S-A1-V1 versus ML are presented in Table 5 and Table 6, respectively. The interaction of age with group is highly significant in both methods, reflecting a steeper average trend in number of exercises with age in the patient group. The fixed intercept and group effects are not significant at a 5% level. Parameters estimates from Algorithm 3S-A1-V1 are

Table 5. Results of algorithm 3S-A1-V1 fitting the LMM (44) to the Blackmore data.

<i>parameters</i>	<i>est.</i>	[2.5% – 97.5%]	<i>bias</i>	<i>std.dev</i>	$\sqrt{\text{MSE}}$	<i>t-stat</i>	<i>p-value</i>
$\beta_{\text{intercept}}$	-0.2754	[-0.6371, 0.09391]	0.0007563	0.1855	0.1855	-1.485	1.376e – 01
β_{age8}	0.06366	[0.005921, 0.1248]	0.001169	0.03069	0.03071	2.073	3.818e – 02
$\beta_{\text{group:patient}}$	-0.3552	[-0.8579, 0.1008]	-0.0008574	0.2406	0.2406	-1.476	1.399e – 01
$\beta_{\text{age8*group:patient}}$	0.2405	[0.1610, 0.3200]	-0.001338	0.03956	0.03958	6.077	1.225e – 09
$\sigma_{\text{intercept}}^2$	2.097	[1.668, 2.651]	0.06293	0.2451	0.2531	8.288	1.155e – 16
σ_{age8}^2	0.02954	[0.02174, 0.05089]	0.006141	0.007498	0.009692	3.048	2.301e – 03
$\text{corr}_{\text{intercept,age8}}$	-0.3006	[-0.5171, -0.1187]	-0.03928	0.3193	0.3306	-2.751	5.936e – 03
σ_{ε}^2	1.533	[1.289, 1.678]	-0.04473	0.09846	0.1081	14.17	1.336e – 45
ρ	2.297	[1.692, 3.490]	0.2573	0.4720	0.5376	4.272	1.937e – 05

$\text{corr}_{\text{interceptage8}}$: correlation coefficient between random effects variables intercept and age8.

Table 6. Results of Gaussian ML fitting of the LMM (44) to the Blackmore data.

<i>parameters</i>	<i>est.</i>	<i>std.dev</i>	<i>t-stat</i>	<i>parameters</i>	<i>est.</i>
$\beta_{\text{intercept}}$	-0.2762	0.1816	-1.521	$\sigma_{\text{intercept}}^2$	2.058
β_{age8}	0.06412	0.03122	2.054	σ_{age8}^2	0.02644
$\beta_{\text{group:patient}}$	-0.3536	0.2343	-1.510	$\text{corr}_{\text{intercept,age8}}$	-0.28
$\beta_{\text{age8*group:patient}}$	0.2396	0.03922	6.110	σ_{ε}^2	1.548

very close to those from ML. The p -value for the random effects ratio estimate is less than 0.001, confirming the presence, w.r.t. the girls, of random effects on the intercept and/or age8, and probably both since the estimates of their variances $\sigma_{\text{intercept}}^2$ and σ_{age8}^2 are significantly greater than zero. It also appears that there is a significant negative correlation between those two random effects.

In the Supplementary material, we fitted another LMM to this data set by adding a random slope for age8*group in (44). The results using Algorithm 3S-A1-V1 show that, in addition to the parameters already significant in (44), the same is true of the variance of the interaction age8*group, implying a probable random effect on it also. We remark that Gaussian ML fitting of this model failed to converge with the default settings in the R function lmer.

7. Concluding Remarks

Till today, it has been difficult to routinely fit LMMs without assuming both random effects and residual errors to have Gaussian distributions (the default in almost all statistical software packages designed for that purpose). Yet, for many data sets, that assumption may be debatable, especially for the random effects. This is disturbing since modeling of random effects behavior is one of the main goals of LMM fitting in the first place. Therefore, there has been an implicit need, for long now, to develop fitting methods for LMMs not requiring Gaussian assumptions, while being applicable to as wide a range of LMMs as possible. Being restricted to variance components models, the venerable ANOVA methods and Rao's MINQUE largely fall short in that respect. In the work presented here, we were able to devise a new iterative fitting methodology for 2-level (or longitudinal) LMMs with only added assumption (to the basic ones) that the residual errors were uncorrelated and homoscedastic. Each variant of that estimation methodology iterates through a small set of estimating equations and, when

convergent, yields nonnegative estimates of variances and SPSS estimates of covariance matrices. Though no Gaussian assumption is involved in the derivation of these EEs, we, however, showed that if the 2-level LMM is, indeed, Gaussian with i.i.d. errors, then these EEs are equivalent to the likelihood ones, safe for a denominator $n-p$ in lieu of n in the EE for the residual variance. Furthermore, the newly developed estimation methodology for LMMs, nicknamed 3S, is not exclusive to 2-level ones. The same ideas can be used to fit some other classes of LMMs as well, as we showed for variance components (or ANOVA) LMMs, generalizing an old method of Henderson for ML estimation in such models of Gaussian type. An interesting by-product we got is also, actually, an extension of that Henderson method to Gaussian ML fitting of 2-level LMMs with u.h.o. errors.

Proof of Theorem 7

In addition to Lemmas S3.1 and S3.2, we shall need the following one:

Lemma 8. *In the 2-level LMM (2), if the clusters random effects matrix \mathbf{D} satisfies:*

$$\mathbf{D} = \frac{1}{m} \sum_{j=1}^m (\mathbf{V}_j^* + \tilde{\mathbf{U}}_j \tilde{\mathbf{U}}_j^T), \tag{45}$$

with $\tilde{\mathbf{U}}_1, \dots, \tilde{\mathbf{U}}_m \in \mathbb{R}^m$, then, letting $\tilde{\mathbf{U}} = (\tilde{\mathbf{U}}_1^T, \dots, \tilde{\mathbf{U}}_m^T)^T \in \mathbb{R}^q$,

$$\text{tr}(\mathbf{Z}\mathbf{G}\mathbf{Z}^T\mathbf{V}^{-1}) = \|\tilde{\mathbf{U}}\|_{\mathbf{G}^{-1}}^2. \tag{46}$$

Proof. From the proof of Proposition S1.4, we already deduce that

$$\text{tr}(\mathbf{Z}\mathbf{G}\mathbf{Z}^T\mathbf{V}^{-1}) = \text{tr}[(\mathbf{G} - \mathbf{V}^*)\mathbf{G}^{-1}]. \tag{47a}$$

Now, in the 2-level LMM (2), we know that $\mathbf{G} = \text{diag}(\mathbf{D}, \dots, \mathbf{D})$, and $\mathbf{V}^* = \text{diag}(\mathbf{V}_1^*, \dots, \mathbf{V}_m^*)$, which are two conformal $q \times q$ block-diagonal matrices. Hence, one also has:

$$\mathbf{G} - \mathbf{V}^* = \text{diag}(\mathbf{D} - \mathbf{V}_1^*, \dots, \mathbf{D} - \mathbf{V}_m^*) \text{ and } \mathbf{G}^{-1} = \text{diag}(\mathbf{D}^{-1}, \dots, \mathbf{D}^{-1}),$$

which implies that $(\mathbf{G} - \mathbf{V}^*)\mathbf{G}^{-1} = \text{diag}((\mathbf{D} - \mathbf{V}_1^*)\mathbf{D}^{-1}, \dots, (\mathbf{D} - \mathbf{V}_m^*)\mathbf{D}^{-1})$. Therefore,

$$\text{tr}[(\mathbf{G} - \mathbf{V}^*)\mathbf{G}^{-1}] = \sum_{j=1}^m \text{tr}[(\mathbf{D} - \mathbf{V}_j^*)\mathbf{D}^{-1}] = \text{tr}\left[\sum_{j=1}^m (\mathbf{D} - \mathbf{V}_j^*) \cdot \mathbf{D}^{-1}\right]. \tag{47b}$$

Notice then that (45) is equivalent to $\sum_{j=1}^m (\mathbf{D} - \mathbf{V}_j^*) = \sum_{j=1}^m \tilde{\mathbf{U}}_j \tilde{\mathbf{U}}_j^T$, which, inserted in (47b), yields:

$$\text{tr}[(\mathbf{G} - \mathbf{V}^*)\mathbf{G}^{-1}] = \text{tr}\left[\sum_{j=1}^m \tilde{\mathbf{U}}_j \tilde{\mathbf{U}}_j^T \cdot \mathbf{D}^{-1}\right] = \sum_{j=1}^m \text{tr}(\tilde{\mathbf{U}}_j \tilde{\mathbf{U}}_j^T \mathbf{D}^{-1}). \tag{47c}$$

Now, for $j = 1, \dots, m$, $\text{tr}(\tilde{\mathbf{U}}_j \tilde{\mathbf{U}}_j^T \mathbf{D}^{-1}) = \text{tr}(\tilde{\mathbf{U}}_j^T \mathbf{D}^{-1} \tilde{\mathbf{U}}_j) = \tilde{\mathbf{U}}_j^T \mathbf{D}^{-1} \tilde{\mathbf{U}}_j \in \mathbb{R}$. Hence, (47c) implies: $\text{tr}[(\mathbf{G} - \mathbf{V}^*)\mathbf{G}^{-1}] = \sum_{j=1}^m \tilde{\mathbf{U}}_j^T \mathbf{D}^{-1} \tilde{\mathbf{U}}_j = \tilde{\mathbf{U}}^T \mathbf{G}^{-1} \tilde{\mathbf{U}} = \|\tilde{\mathbf{U}}\|_{\mathbf{G}^{-1}}^2$, which, with (47a), yields (46). \square

- We now proceed properly to prove Theorem 7:

Proof. Let $\tilde{\beta}$ and \tilde{U} satisfy the HMMEs (10) given σ_ϵ^2 , \mathbf{D} and $\mathbf{Y} = \mathbf{y}$ in the 2-level LMM (2). Then, thanks to Theorem 1,

$$\tilde{\beta} = (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1} \mathbf{y} \text{ and } \tilde{U} = \mathbf{GZ}^T \mathbf{V}^{-1} (\mathbf{y} - \mathbf{X} \tilde{\beta}), \tag{48a}$$

the latter being equivalent here to:

$$\tilde{U}_j = \mathbf{DZ}_j^T \mathbf{V}_j^{-1} (\mathbf{y}_j - \mathbf{X}_j \tilde{\beta}), \text{ for } j = 1, \dots, m. \tag{48b}$$

Now, let Assumptions $\mathcal{A}1 - \mathcal{A}5$ be true. Then, for $j = 1, \dots, m$, $\mathbf{Y}_j \stackrel{\mathcal{L}}{\sim} \mathcal{N}_{n_j}(\mathbf{X}_j \beta, \mathbf{V}_j)$, with density:

$$f_{\mathbf{Y}_j}(\mathbf{y}_j | \beta, \sigma_\epsilon^2, \mathbf{D}) = (2\pi)^{-n_j/2} |\mathbf{V}_j|^{-1/2} \exp\left[-\frac{1}{2} \|\mathbf{y}_j - \mathbf{X}_j \beta\|_{\mathbf{V}_j^{-1}}^2\right],$$

The log-likelihood over the data from all the m clusters is thus:

$$\ell(\beta, \sigma_\epsilon^2, \mathbf{D} | \mathbf{y}) = c - \frac{1}{2} \sum_{j=1}^m \left[\log |\mathbf{V}_j| + (\mathbf{y}_j - \mathbf{X}_j \beta)^T \mathbf{V}_j^{-1} (\mathbf{y}_j - \mathbf{X}_j \beta) \right]. \tag{48c}$$

1) Now, assume also that (33a) holds, hence $\beta = \tilde{\beta}$, thanks to (34). On the other hand, using (S:3.1a) and (S:3.1b) in Lemma S3.2 to differentiate (48c) w.r.t. \mathbf{D} gives:

$$-2 \frac{\partial \ell}{\partial \mathbf{D}} = \sum_{j=1}^m \left[\frac{\partial \log |\mathbf{V}_j|}{\partial \mathbf{D}} + \frac{\partial (\mathbf{y}_j - \mathbf{X}_j \beta)^T \mathbf{V}_j^{-1} (\mathbf{y}_j - \mathbf{X}_j \beta)}{\partial \mathbf{D}} \right] = \text{diag}(S) - 2S,$$

where $S = \sum_{j=1}^m S_j$, with $S_j = W_j - B_j$, $B_j = \mathbf{Z}_j^T \mathbf{V}_j^{-1} \mathbf{Z}_j$, $W_j = \mathbf{Z}_j^T \mathbf{V}_j^{-1} (\mathbf{y}_j - \mathbf{X}_j) (\mathbf{y}_j - \mathbf{X}_j)^T \mathbf{V}_j^{-1} \mathbf{Z}_j$. With Lemma S3.1, the nonsingularity of \mathbf{D} , the fact that $\beta = \tilde{\beta}$ and using (48b), then (S:1.6c),

$$\begin{aligned} \frac{\partial \ell}{\partial \mathbf{D}} = 0 &\Leftrightarrow S = 0 \Leftrightarrow \sum_{j=1}^m W_j = \sum_{j=1}^m B_j \Leftrightarrow \mathbf{D} \cdot \sum_{j=1}^m W_j \cdot \mathbf{D} = \mathbf{D} \cdot \sum_{j=1}^m B_j \cdot \mathbf{D} \\ &\Leftrightarrow \sum_{j=1}^m \tilde{U}_j \tilde{U}_j^T = \sum_{j=1}^m (\mathbf{D} - \mathbf{V}_j^*) \Leftrightarrow \mathbf{D} = \frac{1}{m} \sum_{j=1}^m (\mathbf{V}_j^* + \tilde{U}_j \tilde{U}_j^T) = \tilde{\mathbf{D}}. \end{aligned} \tag{49a}$$

2) Now, let also (33a) and (33b) both hold. Hence $\tilde{\beta} = \beta$ and (49a) are both true. Furthermore, (49a) implies that $\mathbf{G} = \tilde{\mathbf{G}} = \text{diag}(\tilde{\mathbf{D}}, \dots, \tilde{\mathbf{D}})$. Using (S:3.1c) and (S:3.1d) in Lemma S3.2 to differentiate (48c) w.r.t. σ_ϵ^2 ,

$$-2 \frac{\partial \ell}{\partial \sigma_\epsilon^2} = \text{tr}(\mathbf{V}^{-1}) - \|\mathbf{V}^{-1} (\mathbf{y} - \mathbf{X} \beta)\|^2. \tag{49b}$$

Now, given that $\beta = \tilde{\beta}$ and (48a), we get, from (12a) in Lemma 3:

$$\sigma_\epsilon^2 \|\mathbf{V}^{-1} (\mathbf{y} - \mathbf{X} \tilde{\beta})\| = \|\mathbf{y} - \mathbf{X} \tilde{\beta} - \mathbf{Z} \tilde{U}\|. \tag{49c}$$

Also, using (11), then Lemma 8 (given that (49a) is true, by assumption),

$$\sigma_\epsilon^2 \text{tr}(\mathbf{V}^{-1}) = \text{tr}(\mathbf{I}_n) - \text{tr}(\mathbf{ZGZ}^T \mathbf{V}^{-1}) = n - \|\tilde{U}\|_{\tilde{\mathbf{G}}^{-1}}^2. \tag{49d}$$

Inserting (49c) and (49d) in (49b), we get:

$$-2 \frac{\partial \ell}{\partial \sigma_\epsilon^2} = \frac{n - \|\tilde{U}\|_{\tilde{\mathbf{G}}^{-1}}^2}{\sigma_\epsilon^2} - \frac{1}{\sigma_\epsilon^4} \|\mathbf{y} - \mathbf{X} \tilde{\beta} - \mathbf{Z} \tilde{U}\|^2.$$

Consequently, from Corollary 1,

$$\frac{\partial \ell}{\partial \sigma_\varepsilon^2} = 0 \Leftrightarrow n\sigma_\varepsilon^2 = \|\mathbf{y} - \mathbf{X}\tilde{\boldsymbol{\beta}} - \mathbf{Z}\tilde{\mathbf{U}}\|^2 + \sigma_\varepsilon^2 \|\tilde{\mathbf{U}}\|_{\tilde{\mathbf{G}}^{-1}}^2 = \mathbf{y}^T (\mathbf{y} - \mathbf{X}\tilde{\boldsymbol{\beta}} - \mathbf{Z}\tilde{\mathbf{U}}). \quad \square$$

Conflicts of Interest

The authors declare no conflicts of interest regarding the publication of this paper.

References

- [1] Gibbons, R.D. and Hedeker, D. (2000) Applications of Mixed-Effects Models in Biostatistics. *Sankhy*, **62**, 70-103.
- [2] Yarkiner, Z., Hunter, G., O'Neil, R. and de Lusignan, S. (2013) Applications of Mixed Models for Investigating Progression of Chronic Disease in a Longitudinal Dataset of Patient Records from General Practice. *Journal of Biometrics & Biostatistics*, **S9**, Article No. 001. <https://doi.org/10.4172/2155-6180.S9-001>
- [3] Van der Merwe, A.J. and Pretorius, A.L. (2003) An Application of the Mixed Linear Model and the Dirichlet Process Prior in Veterinary Medicine Research. *Journal of Agricultural, Biological, and Environmental Statistics*, **8**, Article No. 328. <https://doi.org/10.1198/1085711032192>
- [4] Milkevych, V., Madsen, P., Gao, H. and Jensen, J. (2021) The Relative Effect of Genomic Information on Efficiency of Bayesian Analysis of the Mixed Linear Model with Unknown Variance. *Journal of Animal Breeding and Genetics*, **138**, 14-22.
- [5] Torkashvand, E., Jozani, M.J. and Torabi, M. (2017) Clustering in Small Area Estimation with Area Level Linear Mixed Models. *Journal of the Royal Statistical Society Series A (Statistics in Society)*, **180**, 1253-1279. <https://doi.org/10.1111/rssa.12308>
- [6] Maiti, T. (2001) Robust Generalized Linear Mixed Models for Small Area Estimation. *Journal of Statistical Planning and Inference*, **98**, 225-238. [https://doi.org/10.1016/S0378-3758\(00\)00302-5](https://doi.org/10.1016/S0378-3758(00)00302-5)
- [7] Kizilkaya, K., Garrick, D.J., Fernando, R.L., Mestav, B. and Yildiz, M.A. (2010) Use of Linear Mixed Models for Genetic Evaluation of Gestation Length and Birth Weight Allowing for Heavy-Tailed Residual Effects. *Genetics Selection Evolution*, **42**, Article No. 26. <https://doi.org/10.1186/1297-9686-42-26>
- [8] Dandine-Roulland, C. and Perdry, H. (2015) The Use of the Linear Mixed Model in Human Genetics. *Human Heredity*, **80**, 196-206. <https://doi.org/10.1159/000447634>
- [9] Lussetti, D., Kuljus, K., Ranneby, B., Ilstedt, U., Falck, J. and Karlsson, A. (2019) Using Linear Mixed Models to Evaluate Stand Level Growth Rates for Dipterocarps and *Macaranga* Species Following Two Selective Logging Methods in Sabah, Borneo. *Forest Ecology and Management*, **437**, 372-379. <https://doi.org/10.1016/j.foreco.2019.01.044>
- [10] Zamudio, F., Yanez, M., Guerra, F., Fuentes, D. and Gonzalez, A. (2020) Comparative Analysis of SNP Data and Hybrid Taxa Information by Using a Classificatory Linear Mixed Model to Study the Genetic Variation and Heritability of Initial Height Growth in Selected Poplar Hybrids. *Tree Genetics & Genomes* **16**, Article No. 69. <https://doi.org/10.1007/s11295-020-01435-1>
- [11] Jiang, J. (2007) Linear and Generalized Linear Mixed Models and Their Applications. Springer, New York.
- [12] Demidenko, E. (2013) Mixed Models: Theory and Applications with R. 2nd Edition,

John Wiley & Sons, Inc., Hoboken.

- [13] West, B.T., Welch, K.B. and Galecki, A.T. (2007) *Linear Mixed Models: A Practical Guide Using Statistical Software*. Chapman & Hall/CRC, London.
- [14] Searle, S.R., Casella, G. and McCulloch, C.E. (2006) *Variance Components*. Wiley, New York.
- [15] Hartley, H.O. and Rao, J.N.K. (1967) Maximum Likelihood Estimation for the Mixed Analysis of Variance Model. *Biometrika*, **54**, 93-108. <https://doi.org/10.1093/biomet/54.1-2.93>
- [16] Harville, D.A. (1977) Maximum Likelihood Approaches to Variance Component Estimation and to Related Problems. *Journal of the American Statistical Association*, **72**, 320-338. <https://doi.org/10.1080/01621459.1977.10480998>
- [17] Lindstrom, M.J. and Bates, D.M. (1988) Newton-Raphson and EM Algorithms for Linear Mixed-Effects Models for Repeated-Measures Data. *Journal of the American Statistical Association*, **83**, 1014-1022.
- [18] Pinheiro, J.C. and Bates, D.M. (2000) *Mixed-Effects Models in S and S-PLUS*. Springer, NY.
- [19] Gumedze, F.N. and Dunne, T.T. (2011) Parameter Estimation and Inference in the Linear Mixed Model. *Linear Algebra and its Applications*, **435**, 1920-1944. <https://doi.org/10.1016/j.laa.2011.04.015>
- [20] Dempster, A.P., Laird, N.M. and Rubin, D.B. (1977) Maximum Likelihood from Incomplete Data via the EM Algorithm. *Journal of the Royal Statistical Society Series B*, **39**, 1-38.
- [21] Laird, N.M. (1982) Computation of Variance Components Using the EM Algorithm. *Journal of Statistical Computation and Simulation*, **14**, 295-303. <https://doi.org/10.1080/00949658208810550>
- [22] Laird, N.M. and Ware, J.M. (1982) Random Effects Models for Longitudinal Data. *Biometrics*, **38**, 963-974. <https://doi.org/10.2307/2529876>
- [23] Lange, N. and Ryan, L. (1989) Assessing Normality in Random Effects Models. *Annals of Statistics*, **17**, 624-642. <https://doi.org/10.1214/aos/1176347130>
- [24] Li, S., Cai, T.T. and Li, H. (2021) Inference for High-Dimensional Linear Mixed-Effects Models: A Quasi-Likelihood Approach. *Journal of the American Statistical Association*. (In Press) <https://doi.org/10.1080/01621459.2021.1888740>
- [25] Eisenhart, C. (1947) The Assumptions Underlying the Analysis of Variance. *Biometrics*, **3**, 1-21. <https://doi.org/10.2307/3001534>
- [26] Anderson, R.L. and Bancroft, T.A. (1952) *Statistical Theory in Research*. McGraw-Hill, New York.
- [27] Montgomery, D.C. (2004) *Design and Analysis of Experiments*. Wiley, New York.
- [28] Oehlert, G.W. (2010) *A First Course in Design and Analysis of Experiments*. Library of Congress Cataloging-in-Publication Data, USA.
- [29] Henderson, C.R. (1953) Estimation of Variance and Covariance Components. *Biometrics*, **9**, 226-252. <https://doi.org/10.2307/3001853>
- [30] Rao, C.R. (1970) Estimation of Heteroscedastic Variances in Linear Models. *Journal of the American Statistical Association*, **65**, 161-172. <https://doi.org/10.1080/01621459.1970.10481070>
- [31] Rao, C.R. (1971) Estimation of Variance and Covariance Components: MINQUE Theory. *Journal of Multivariate Analysis*, **1**, 257-275. [https://doi.org/10.1016/0047-259X\(71\)90001-7](https://doi.org/10.1016/0047-259X(71)90001-7)

- [32] Rao, C.R. (1971) Minimum Variance Quadratic Unbiased Estimation of Variance Components. *Journal of Multivariate Analysis*, **1**, 445-456. [https://doi.org/10.1016/0047-259X\(71\)90019-4](https://doi.org/10.1016/0047-259X(71)90019-4)
- [33] Rao, C.R. (1972) Estimation of Variance and Covariance Components in Linear Models. *Journal of the American Statistical Association*, **67**, 112-115. <https://doi.org/10.1080/01621459.1972.10481212>
- [34] Rao, C.R. and Kleffe, J. (1988) Estimation of Variance Components and Applications. North-Holland, Amsterdam.
- [35] Jiang, J. (1996) REML Estimation: Asymptotic Behavior and Related Topics. *Annals of Statistics*, **24**, 255-286. <https://doi.org/10.1214/aos/1033066209>
- [36] Heyde, C.C. (1994) A Quasi-Likelihood Approach to the REML Estimating Equations. *Statistics & Probability Letters*, **21**, 381-384. [https://doi.org/10.1016/0167-7152\(94\)00035-2](https://doi.org/10.1016/0167-7152(94)00035-2)
- [37] Jiang, J., Luan, Y. and Wang, Y. (2007) Iterative Estimating Equations: Linear Convergence and Asymptotic Properties. *The Annals of Statistics*, **35**, 2233-2260. <https://doi.org/10.1214/009053607000000208>
- [38] Goldstein, H. (1986) Multilevel Mixed Linear Model Analysis Using Iterative Generalized Least Squares. *Biometrika*, **73**, 43-56. <https://doi.org/10.1093/biomet/73.1.43>
- [39] R Core Team (2019) R: A Language and Environment for Statistical Computing. R Foundation for Statistical Computing, Vienna. <https://www.R-project.org/>
- [40] Searle, S.R. (1995) The Matrix Handling of BLUE and BLUP in the Mixed Linear Model. BU-1275-MA, Biometrics Unit, Cornell University, Ithaca.
- [41] Henderson, C.R. (1950) Estimation of Genetic Parameters (Abstract). *Annals of Mathematical Statistics*, **21**, 309-310.
- [42] Henderson, C.R., Kempthorne, O., Searle, S.R. and von Krosigk, C.N. (1959) The Estimation of Environmental and Genetic Trends from Records Subject to Culling. *Biometrics*, **15**, 192-218. <https://doi.org/10.2307/2527669>
- [43] Henderson, C.R. (1973) Sire Evaluation and Genetic Trends. *Journal of Animal Science*, 10-41. <https://doi.org/10.1093/ansci/1973.Symposium.10>
- [44] Henderson, C.R. (1963) Selection Index and Expected Genetic Advance. In: *Statistical Genetics and Plant Breeding, National Academy of Sciences*, No. 982, National Research Council Publication, Washington DC, 141-163.
- [45] Carpenter, J.R., Goldstein, H. and Rasbash, J. (2003) A Novel Bootstrap Procedure for Assessing the Relationship between Class Size and Achievement. *Journal of the Royal Statistical Society. Series C*, **52**, 431-443. <https://doi.org/10.1111/1467-9876.00415>
- [46] Van der Leeden, R., Meijer, E. and Busing, F. (2008) Chapter 11. Resampling Multilevel Models. In: de Leeuw, J. and Meijer, E., Eds., *Handbook of Multilevel Analysis*, Springer, New York, 401-433.
- [47] Chambers, R. and Chandra, H. (2013) A Random Effect Block Bootstrap for Clustered Data. *Journal of Computational and Graphical Statistics*, **22**, 452-470. <https://doi.org/10.1080/10618600.2012.681216>
- [48] Thai, H.-T., Mentré, F., Holford, N., Veyrat-Follet, C. and Comets, E. (2013) A Comparison of Bootstrap Approaches for Estimating Uncertainty of Parameters in Linear Mixed-Effects Models. *Pharmaceutical Statistics*, **12**, 129-140. <https://doi.org/10.1002/pst.1561>
- [49] Modugno, L. and Giannerini, T. (2015) The Wild Bootstrap for Multilevel Models. *Communications in Statistics— Theory and Methods*, **44**, 4812-4825.

- <https://doi.org/10.1080/03610926.2013.802807>
- [50] Bates, D., Maechler, M., Bolker, B. and Walker, S. (2015) Fitting Linear Mixed-Effects Models Using lme4. *Journal of Statistical Software*, **67**, 1-48.
<https://doi.org/10.18637/jss.v067.i01>
- [51] Cook, F.E. (1938) Chocolate Cake, I. Optimum Baking Temperature. Master's Thesis, Iowa State College, Iowa.
- [52] Efron, B. and Tibsharani, J. (1993) An Introduction to the Bootstrap. Chapman and Hall, London.
- [53] Davis, C., Blackmore, E., Katzman, D.K. and Fox, J. (2005) Female Adolescents with Anorexia Nervosa and Their Parents: A Case-Control Study of Exercise Attitudes and Behaviours. *Psychological Medicine*, **35**, 377-386.
<https://doi.org/10.1017/S0033291704003447>

Supplementary Material

Supplementary Information for the article: Nonparametric Estimation in Linear Mixed Models with Uncorrelated Homoscedastic Errors

Introduction

In this document, we present supplementary materials useful for the understanding of the article. For that, the first section gathers some known results on LMMs scattered here and there in the literature, and some new results of our own. Section S2 presents some important implementation details about our presented iterative methods for fitting LMMs, while Section S3 presents two lemmas useful for proving Theorem 7 in the article. In Section S4, 3S fitting algorithms for 2-level LMMs with u.h.o. errors and diagonal covariance matrix for the cluster random effects, and ANOVA LMMs are presented. More analysis of the Blackmore data is presented in the last section.

S1. LMMs without Gaussian Assumptions: Some Useful Results

S1.1. An Equivalent Formulation of Assumptions $\mathcal{A}1$ - $\mathcal{A}3$

It is useful to note the following equivalent formulation of Assumptions $\mathcal{A}1$ - $\mathcal{A}3$:

Proposition S1.1. *Assumptions $\mathcal{A}1$ - $\mathcal{A}3$ are equivalent to the following set of 3 assumptions:*

- a1) *The U_j 's are identically distributed in \mathbb{R}^r .*
- a2) *$\forall j$, the two random vectors U_j and ε_j are independent.*
- a3) *The random couples $(U_1, \varepsilon_1), \dots, (U_m, \varepsilon_m)$ are mutually independent.*

S1.2. More about the BLP of $U | Y$ in an LMM

The best linear predictor $\bar{U} = \text{BLP}(U | Y)$, given by (5), plays a key role in our new estimation methodology for LMMs with u.h.o. errors. That is why we expand on it here. And, again, no Gaussian distributional assumption is involved.

S1.2.1. Covariance Matrix and Another Expression of the BLP of $U | Y$ in an LMM

The relationship between the covariance matrix of $\bar{U} = \text{BLP}(U | Y)$ and that of the true random effects vector U will be instrumental in estimating the latter in our methodology, starting with:

$$\mathbb{M} \text{cov}(\bar{U}) = \mathbf{GZ}^T \mathbf{V}^{-1} \mathbb{M} \text{cov}(Y) \mathbf{V}^{-1} \mathbf{ZG} = \mathbf{GZ}^T \mathbf{V}^{-1} \mathbf{V} \mathbf{V}^{-1} \mathbf{ZG} = \mathbf{GZ}^T \mathbf{V}^{-1} \mathbf{ZG}, \quad (\text{S:1.1a})$$

thanks to (5). Hence,

$$\mathbb{M} \text{cov}(\bar{U}) = \mathbb{E}(\bar{U} \bar{U}^T) = \mathbf{G} - \mathbf{V}^*, \quad (\text{S:1.1b})$$

where \mathbf{V}^* is the symmetric matrix:

$$\mathbf{V}^* = \mathbf{G} - \mathbf{GZ}^T \mathbf{V}^{-1} \mathbf{ZG}. \quad (\text{S:1.1c})$$

Introducing the apparently complicated matrix V^* is motivated by the fact that it is only $q \times q$ (and not $n \times n$ like V) and turns out instrumental in devising an alternative useful way of rewriting the BLP (5) in the LMM (1), more so with u.h.o. residual errors:

Lemma S1.2. *In the LMM (1) with Assumption $A1_g$, the following identities hold:*

$$V^* = (Z^T R^{-1} Z + G^{-1})^{-1}, \tag{S:1.2a}$$

$$GZ^T V^{-1} = V^* Z^T R^{-1}. \tag{S:1.2b}$$

Moreover, if Assumption $A4_g$ also holds, then these identities simplify to:

$$V^* = (\sigma_\varepsilon^{-2} Z^T Z + G^{-1})^{-1} = \sigma_\varepsilon^2 \cdot (Z^T Z + \sigma_\varepsilon^2 G^{-1})^{-1}, \tag{S:1.2c}$$

$$GZ^T V^{-1} = \sigma_\varepsilon^{-2} V^* Z^T. \tag{S:1.2d}$$

The advantage of (S:1.2c)-(S:1.2d) is that they give respective expressions of V^* and $GZ^T V^{-1}$ using only $q \times q$ and $n \times q$ matrices, and no $n \times n$ one. These imply a much computationally cheaper formula for the BLP of $U | Y$ in an LMM with u.h.o. errors, when compared with (5):

Proposition S1.3. *In the LMM (1) with Assumption $A1_g$, one has:*

$$\bar{U} = \text{BLP}(U | Y) = V^* Z^T R^{-1} (Y - X \beta). \tag{S:1.3a}$$

Moreover, if Assumption $A4_g$ also holds, then this simplifies to:

$$\bar{U} = \text{BLP}(U | Y) = \sigma_\varepsilon^{-2} V^* Z^T (Y - X \beta). \tag{S:1.3b}$$

It is also worth noting:

Proposition S1.4. *In the LMM (1) with Assumption $A1_g$, one has:*

$$\mathbb{E}(\|\bar{U}\|_{G^{-1}}^2) = q - \text{tr}(G^{-1} V^*). \tag{S:1.4a}$$

Moreover, if Assumption $A4_g$ also holds, then we also have:

$$\mathbb{E}(\|\bar{U}\|_{G^{-1}}^2) = n - \sigma_\varepsilon^2 \text{tr}(V^{-1}) = \sigma_\varepsilon^{-2} \text{tr}(V^* Z^T Z). \tag{S:1.4b}$$

Proof. Given that $\|\bar{U}\|_{G^{-1}}^2 = \bar{U}^T G^{-1} \bar{U}$, $\mathbb{E}(\bar{U}) = \mathbf{0}$ and (S:1.1b),

$$\begin{aligned} \mathbb{E}(\|\bar{U}\|_{G^{-1}}^2) &= \mathbb{E}(\bar{U})^T G^{-1} \mathbb{E}(\bar{U}) + \text{tr}[G^{-1} \mathbb{M} \text{cov}(\bar{U})] = \text{tr}[G^{-1} \mathbb{M} \text{cov}(\bar{U})] \\ &= \text{tr}[G^{-1} (G - V^*)] = \text{tr}(I_q - G^{-1} V^*) = q - \text{tr}(G^{-1} V^*). \end{aligned}$$

Secondly, from (S:1.1a), $\text{tr}[G^{-1} \mathbb{M} \text{cov}(\bar{U})] = \text{tr}(Z^T V^{-1} Z G) = \text{tr}(V^{-1} Z G Z^T)$.

Now, if Assumption $A4_g$ also holds, then

$$\text{tr}(V^{-1} Z G Z^T) = \text{tr}[V^{-1} (V - \sigma_\varepsilon^2 I_n)] = \text{tr}(I_n) - \sigma_\varepsilon^2 \text{tr}(V^{-1}),$$

which entails the first equality in (S:1.4b). For the second one, using (S:1.2d),

$$\text{tr}(V^{-1} Z G Z^T) = \text{tr}(Z G Z^T V^{-1}) = \sigma_\varepsilon^{-2} \text{tr}(Z V^* Z^T) = \sigma_\varepsilon^{-2} \text{tr}(V^* Z^T Z).$$

The last proposition implies some additional striking matrix identities in an LMM with u.h.o. errors:

Proposition S1.5. In the LMM(1) with Assumptions $\mathcal{A}1_g$ and $\mathcal{A}4_g$,

$$\sigma_\epsilon^2 \text{tr}(\mathbf{V}^{-1}) + \sigma_\epsilon^{-2} \text{tr}(\mathbf{V}^* \mathbf{Z}^T \mathbf{Z}) = n, \tag{S:1.5a}$$

$$\text{tr}(\mathbf{G}^{-1} \mathbf{V}^*) + \sigma_\epsilon^{-2} \text{tr}(\mathbf{V}^* \mathbf{Z}^T \mathbf{Z}) = q. \tag{S:1.5b}$$

S1.2.2. Cluster Partitioning of the BLP($U | Y$) in a 2-Level LMM

In the 2-level LMM (2), we can partition $\bar{U} = \text{BLP}(U | Y)$ cluster-wise:

$$\bar{U} = \begin{pmatrix} \bar{U}_1 \\ \vdots \\ \bar{U}_m \end{pmatrix}, \text{ with each } \bar{U}_j \text{ a random } r\text{-vector } (j = 1, \dots, m). \tag{S:1.6a}$$

Combining that with (5), the assumptions and the partitioned structures of the various design and covariance matrices \mathbf{G} , \mathbf{R} and \mathbf{V} , we get:

$$\bar{U}_j = \text{BLP}(U_j | Y_j) = \mathbf{DZ}_j^T \mathbf{V}_j^{-1} (\mathbf{Y}_j - \mathbf{X}_j \beta) \quad (j = 1, \dots, m). \tag{S:1.6b}$$

Similarly, from (S:1.1c), one sees that $\mathbf{V}^* = \text{diag}(\mathbf{V}_1^*, \dots, \mathbf{V}_m^*)$, where

$$\mathbf{V}_j^* = \mathbf{D} - \mathbf{DZ}_j^T \mathbf{V}_j^{-1} \mathbf{Z}_j \mathbf{D} \in \mathcal{M}_r(\mathbb{R}) \quad (j = 1, \dots, m). \tag{S:1.6c}$$

Furthermore, from (S:1.6b), calculations similar to (S:1.1a) give:

$$\mathbb{M} \text{cov}(\bar{U}_j) = \mathbb{E}(\bar{U}_j \bar{U}_j^T) = \mathbf{DZ}_j^T \mathbf{V}_j^{-1} \mathbf{Z}_j \mathbf{D} = \mathbf{D} - \mathbf{V}_j^* \quad (j = 1, \dots, m). \tag{S:1.6d}$$

So, here, $\mathbb{M} \text{cov}(\bar{U}) = \text{diag}(\mathbb{M} \text{cov}(\bar{U}_1), \dots, \mathbb{M} \text{cov}(\bar{U}_m))$, thus yielding

$$\mathbb{M} \text{cov}(\bar{U}) = \text{diag}(\mathbf{D} - \mathbf{V}_1^*, \dots, \mathbf{D} - \mathbf{V}_m^*) = \mathbf{G} - \mathbf{V}^*. \tag{S:1.6e}$$

Finally, the cluster-wise versions of Lemma S1.2 and Propositions S1.3-S1.5 for the 2-level LMM (2) are:

Lemma S1.6. In the 2-level LMM (2), the following identities hold, for $j = 1, \dots, m$:

$$\mathbf{V}_j^* = (\mathbf{Z}_j^T \mathbf{R}_j^{-1} \mathbf{Z}_j + \mathbf{D}^{-1})^{-1}, \tag{S:1.7a}$$

$$\mathbf{DZ}_j^T \mathbf{V}_j^{-1} = \mathbf{V}_j^* \mathbf{Z}_j^T \mathbf{R}_j^{-1}. \tag{S:1.7b}$$

Moreover, if Assumption $\mathcal{A}4_g$ also holds, then these identities simplify to:

$$\mathbf{V}_j^* = (\sigma_\epsilon^{-2} \mathbf{Z}_j^T \mathbf{Z}_j + \mathbf{D}^{-1})^{-1} = \sigma_\epsilon^2 (\mathbf{Z}_j^T \mathbf{Z}_j + \sigma_\epsilon^2 \mathbf{D}^{-1})^{-1}, \tag{S:1.7c}$$

$$\mathbf{DZ}_j^T \mathbf{V}_j^{-1} = \sigma_\epsilon^{-2} \mathbf{V}_j^* \mathbf{Z}_j^T. \tag{S:1.7d}$$

Proposition S1.7. In the 2-level LMM(2), one has, for $j = 1, \dots, m$:

$$\bar{U}_j = \text{BLP}(U_j | Y_j) = \mathbf{V}_j^* \mathbf{Z}_j^T \mathbf{R}_j^{-1} (\mathbf{Y}_j - \mathbf{X}_j \beta). \tag{S:1.8a}$$

Moreover, if Assumption $\mathcal{A}4_g$ also holds, then this simplifies to:

$$\bar{U}_j = \text{BLP}(U_j | Y_j) = \sigma_\epsilon^{-2} \mathbf{V}_j^* \mathbf{Z}_j^T (\mathbf{Y}_j - \mathbf{X}_j \beta). \tag{S:1.8b}$$

Proposition S1.8. In the 2-level LMM (2), one has, for $j = 1, \dots, m$:

$$\mathbb{E}(\|\bar{U}_j\|_{\mathbf{D}^{-1}}^2) = r - \text{tr}(\mathbf{D}^{-1} \mathbf{V}_j^*). \tag{S:1.9}$$

Moreover, if Assumption $\mathcal{A}4_g$ also holds, then we also have:

$$\mathbb{E}\left(\|\bar{U}_j\|_{D^{-1}}^2\right) = n_j - \sigma_\varepsilon^2 \text{tr}(\mathbf{V}_j^{-1}) = \sigma_\varepsilon^{-2} \text{tr}(\mathbf{V}_j^* \mathbf{Z}_j^T \mathbf{Z}_j). \quad (\text{S:1.10})$$

Proposition S1.9. In the 2-level LMM (2) with Assumptions $\mathcal{A}1_g$ and $\mathcal{A}4_g$, for $j = 1, \dots, m$:

$$\sigma_\varepsilon^2 \text{tr}(\mathbf{V}_j^{-1}) + \sigma_\varepsilon^{-2} \text{tr}(\mathbf{V}_j^* \mathbf{Z}_j^T \mathbf{Z}_j) = n_j, \quad (\text{S:1.11a})$$

$$\text{tr}(\mathbf{D}^{-1} \mathbf{V}_j^*) + \sigma_\varepsilon^{-2} \text{tr}(\mathbf{V}_j^* \mathbf{Z}_j^T \mathbf{Z}_j) = r. \quad (\text{S:1.11b})$$

S2. 3S Iterative Algorithms: Some Important Implementation Details

We give here some useful precisions to effectively program Algorithms 3S-A1-V1 and 3S-A1-V2.

S2.1. Solving the HMMs in Algorithm 3S-A1-V1

Our way of solving a system of the form (29) is based on the following result which we admit:

Theorem S2.1. (Solving a 2×2 block nonsingular system) Let

$A = \begin{pmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{pmatrix} \in \mathcal{M}_n(\mathbb{R})$ be nonsingular, with $n = n_1 + n_2$ and $n_1, n_2 \in \mathbb{N}^*$. If

$A_{11} \in \mathcal{M}_{n_1}(\mathbb{R})$ and is nonsingular, then:

- 1) $B_{22} = A_{22} - A_{21} A_{11}^{-1} A_{12} \in \mathcal{M}_{n_2}(\mathbb{R})$ and is nonsingular.
- 2) The unique solution of the linear system $AX = Y$ in \mathbb{R}^n is given by:

$$\begin{cases} X_2 = B_{22}^{-1} (Y_2 - A_{21} A_{11}^{-1} Y_1) \in \mathbb{R}^{n_2} \\ X_1 = A_{11}^{-1} (Y_1 - A_{12} X_2) \in \mathbb{R}^{n_1}, \end{cases}$$

with $X_1, Y_1 \in \mathbb{R}^{n_1}$, $X_2, Y_2 \in \mathbb{R}^{n_2}$ such that $X = \begin{pmatrix} X_1 \\ X_2 \end{pmatrix}$ and $Y = \begin{pmatrix} Y_1 \\ Y_2 \end{pmatrix}$.

Obviously, Theorem S2.1 has a twin version whereby it is rather the lower diagonal block $A_{22} \in \mathcal{M}_{n_2}(\mathbb{R})$ which is assumed nonsingular. To solve the linear system (29) at the beginning of each new iteration $t+1$ in Algorithm 3S-A1-V1, we use Theorem S2.1 with:

$$\begin{aligned} A_{11} &= \mathbf{X}^T \mathbf{X} \in \mathcal{M}_p(\mathbb{R}), \quad A_{12} = \mathbf{X}^T \mathbf{Z} \in \mathcal{M}_{p,q}(\mathbb{R}), \\ A_{21} &= \mathbf{Z}^T \mathbf{X} = A_{12}^T \in \mathcal{M}_{q,p}(\mathbb{R}), \quad A_{22} = \mathbf{Z}^T \mathbf{Z} + \hat{\sigma}_\varepsilon^{2(t)} \hat{\mathbf{G}}^{(t)-1} \in \mathcal{M}_q(\mathbb{R}), \\ Y_1 &= \mathbf{X}^T \mathbf{y} \in \mathbb{R}^p, \quad Y_2 = \mathbf{Z}^T \mathbf{y} \in \mathbb{R}^q. \end{aligned}$$

Remark S2.1. The product matrices A_{11} , A_{12} , A_{21} , $\mathbf{Z}^T \mathbf{Z}$ and vectors Y_1 , Y_2 are all independent from the iteration index t . Hence, they can be computed before the iterations start, which is a major computational advantage for the running time of Algorithm 3S-A1-V1.

S2.2. Choice of the Matrix Norm $\|\cdot\|^{(M)}$ in Instruction 3 of the 3S Algorithms

In the two iterative estimating algorithms for 2-level LMMs presented in Section

4.2.3, all the matrix iterates $\hat{\mathbf{D}}^{(t)}$ are at least SPSD. Based on that, we choose an appropriate matrix norm $\|\cdot\|^{(M)}$ to use in *Instruction 3* of these algorithms to cheaply compute $\|\hat{\mathbf{D}}^{(t+1)}\|^{(M)}$ at any new iteration $t+1$. Indeed, if $A = (a_{ij}) \in \mathcal{M}_r(\mathbb{R})$ is SPSD, then it is known that:

$$\max_{k,l} |a_{kl}| = \max_k a_{kk}.$$

Consequently, if we take as matrix norm:

$$\|A\|^{(M)} = \max_{k,l} |a_{kl}|, \tag{S:2.1}$$

then $\|A\|^{(M)} = \max[\text{diag}(A)]$, where $\text{diag}(A)$ is the vector of diagonal elements in matrix A . So using the matrix norm given by (S:2.1), $\|\hat{\mathbf{D}}^{(t+1)}\|^{(M)}$ is very cheap to compute and the first control test inequality of the stopping criterion in *Instruction 3* then simplifies to:

$$\max_{k,l} \left| \left(\hat{\mathbf{D}}^{(t+1)} - \hat{\mathbf{D}}^{(t)} \right)_{kl} \right| \leq \delta_1 \cdot \max[\text{diag}(\hat{\mathbf{D}}^{(t+1)})].$$

S2.3. Terminating a 3S Iterative Algorithm: *The Various Scenarios*

An iterative algorithm aimed at equations solving (or optimization), however theoretically well crafted, almost always runs the risk of *premature termination* in some cases, i.e. the algorithm stalls before its expected convergence has been achieved. Algorithms 3S-A1-V1 and 3S-A1-V2 are not immune to that threat: each of them may fail on a given data set because the latter is far from satisfying one of the LMM assumptions upon which the algorithm is built. Therefore, it is necessary to include, in the inner workings of our 3S algorithms, capabilities to detect, as early as possible, the most plausible scenarios of abnormal termination or behavior.

So, as is traditional in equations solving and optimization iterative algorithms, we use a *categorical flag variable* code.Stop in each of our algorithms to record the cause of the algorithm termination during a particular run, with success code:

code.Stop = OK:1 *convergence, with satisfaction of the intended stopping criterion.*

Then we handle the following cases of possible premature stoppage, in which we also include a code. Warning flag variable giving rather a specific warning about the iterations (even if they ended apparently well).

S2.3.1. code.Stop = KO:X: *Fixed Effects Design Matrix X Not of Full Column Rank*

As developed till today, a key assumption used in our 3S methodology for estimation in LMMs is $\mathcal{A}2_g$, i.e. the fixed effects design matrix \mathbf{X} is of full column rank. If not, then some of our derivations no longer hold. Although the methodology might be extended in the future to bypass the need of that assumption, at present we stop any of our algorithms if $\mathcal{A}2_g$ is not satisfied at the outset.

1) Checking if the fixed effects design matrix X is full column rank in Algorithm 3S-A1-V1

For Algorithm 3S-A1-V1, the checking is based on the fact that computing and storing the $p \times p$ product matrix $X^T X$ is a necessary requirement in that algorithm, to form the coefficient matrix of the linear system of HMMEs (29). Now, X and $X^T X$ always have the same rank, so X has full column rank if, and only if, $X^T X$ is nonsingular. But the latter holds if, and only if, the symmetric matrix $X^T X$ (which is at least SPSPD) is SPD, which is equivalent to having a *Cholesky factorization*, i.e. $X^T X = U^T U$, where U is a square upper triangular matrix with positive diagonal elements.

Hence, at the beginning of Algorithm 3S-A1-V1, to check for Assumption $\mathcal{A}2_g$, once $X^T X$ has been obtained, we compute its *pivoted Cholesky factorization*, using an available software program. Modern numerical software such as the LAPACK routines (Anderson et al., 1999), called by the R statistical software (our programming tool) functions handling Numerical Linear Algebra (NLA) calculations, can compute that factorization and check whether an input symmetric matrix is even SPSPD in the first place, and if yes, then check if it is likely, or not, SPD, at least up to the numerical instabilities due to the computer rounding errors. They can also output an estimate of the rank k of an SPSPD matrix and indicate its k columns most likely linearly uncorrelated. Applying that to $X^T X$, we can decide whether it is clearly SPD, or likely close to a singular SPSPD matrix.

In the latter case, we set `code.Stop = KO:X`, and do not launch Algorithm 3S-A1-V1 at all, informing the user why. But we also output the estimated rank k of $X^T X$ and the indices of its k columns most likely uncorrelated. Since these two pieces of information are the same for X , the user can then, if he/she wishes, re-launch the 3S iterative algorithm, but with that fixed effects design matrix changed by suppressing its other $p - k$ columns.

2) Checking if the fixed effects design matrix X is full column rank in Algorithm 3S-A1-V2

Algorithm 3S-A1-V2, instead of the HMMEs, needs to solve, at each iteration, an OLS problem which design matrix is X . To achieve that, rather than first computing the product matrix $X^T X$, numerical analysts recommend to compute the QR factorization (Demmel, 1997) of X . Interestingly, modern numerical software such as the LINPACK routines (Dongarra et al., 1978) called by the R NLA functions can do so while also giving an estimate of the rank of X . So for Algorithm 3S-A1-V2, our checking of Assumption $\mathcal{A}2_g$ will be based on the results of the *QR factorization with pivoting* of X , with the conclusion drawn as for Algorithm 3S-A1-V1 above.

S2.3.2. A Vital Issue for LMM Fitting: *Random Effects or Not?*

People trying to fit an LMM rather than an LM to their data most often do not worry whether the incurred excess in mathematical modeling and computational time is worth the price in the first place. They do not ask themselves the obvious

question:

Are there really random effects in the data with the provided clustering?

But that's a mistake, noticeably because it turns out that the answer to that question directly impacts the way an LMM fitting iterative procedure really behaves on a given data set. In the programming of a 3S algorithm, we anticipate three possible scenarios related to that issue.

3) code.Stop = KO:1: RE covariance matrix likely (or close to) singular

Even when the cluster random effects are there, they may have been poorly parameterized, so that their scalar components are correlated, resulting in a rank deficient (or singular) covariance matrix \mathbf{D} . If so, this will affect adversely any of our 3S algorithms since they all assume \mathbf{D} to be SPD.

To deal with that potential difficulty in a 3S algorithm, since \mathbf{D} is not known (we are trying to estimate it among other parameters...), our strategy is that after computing each new iterate $\hat{\mathbf{D}}^{(t)}$ ($t \geq 1$), we test whether it is of full rank or not. Since in our computations, $\hat{\mathbf{D}}^{(t)}$ is always at least SPSD, we numerically check for its rank by attempting a pivoted Cholesky factorization on it along the same lines as described above for the matrix $\mathbf{X}^T \mathbf{X}$ in the code.Stop = KO:X scenario. If the Cholesky software routine declares the newly computed $\hat{\mathbf{D}}^{(t)}$ at an iteration t of not being full rank, the iterative algorithm is stopped, informing the user why: *the true cluster random effects covariance matrix \mathbf{D} is either singular or close to such a matrix.*

After that premature termination (and as for the case of singular $\mathbf{X}^T \mathbf{X}$ previously), the algorithm also returns the estimated rank ℓ of $\hat{\mathbf{D}}^{(t)}$ and the indices of its ℓ most uncorrelated columns. Using these pieces of information as estimates of the corresponding features of the true (but unknown) \mathbf{D} , the user can re-launch the 3S iterative algorithm, but with cluster random effects vectors U_1, \dots, U_m of dimension reduced from r to ℓ , only keeping the ℓ scalar components associated to the ℓ columns most uncorrelated in \mathbf{D} . In the data, this amounts to keeping just the corresponding ℓ columns in each of the cluster random effects design matrices $\mathbf{Z}_1, \dots, \mathbf{Z}_m$.

Note, however, that this does not work when $r = 1$, i.e. each cluster random effect is actually a scalar rather than a vector because, then, the cluster random effects covariance matrix \mathbf{D} also reduces to a (nonnegative) scalar and is, therefore, rank deficient *if, and only if*, that scalar is zero. This will also be true of its $\hat{\mathbf{D}}^{(t)}$ iterates. But the latter matrices will unlikely be numerically close to zero up to computer rounding errors, so there is little chance that any of them would be detected as rank deficient based on the Choleski Factorization. So we need to handle other abnormal possible termination scenarios including that subcase.

4) code.Stop = KO:2: Slow convergence, likely because of no cluster RE

Recall that a 3S algorithm has converged when the 3 inequalities

$$\|\hat{\mathbf{D}}^{(t+1)} - \hat{\mathbf{D}}^{(t)}\|^{(M)} \leq \delta_1 \cdot \|\hat{\mathbf{D}}^{(t+1)}\|^{(M)}, \quad (\text{S:2.2a})$$

$$\|\hat{\beta}^{(t+1)} - \hat{\beta}^{(t)}\| \leq \delta_2 \cdot \|\hat{\beta}^{(t+1)}\|, \quad (\text{S:2.2b})$$

$$\left| \hat{\sigma}_\varepsilon^{2(t+1)} - \hat{\sigma}_\varepsilon^{2(t)} \right| \leq \delta_3 \cdot \hat{\sigma}_\varepsilon^{2(t+1)}, \quad (\text{S:2.2c})$$

are simultaneously satisfied, where $\delta_1, \delta_2, \delta_3$ are relative tolerance levels set in $(0, 1)$, and $\|\cdot\|^{(M)}$ is a chosen matrix norm. Now, when a 3S algorithm appears to be converging slowly, an inspection of the evolution of the iterations intermediate results reveals that this is, more often than not, entirely due to the slow convergence of the sequence of $\hat{\mathbf{D}}^{(t)}$ iterates. This manifests itself through the difficulty of having the relative error control inequality (S:2.2a) satisfied, whereas there is really a decrease in the sequence of absolute errors estimates $\|\hat{\mathbf{D}}^{(t+1)} - \hat{\mathbf{D}}^{(t)}\|^{(M)}$. From Numerical Analysis, we know that such a situation typically happens when the sought limit value of the $\hat{\mathbf{D}}^{(t)}$ iterates (*i.e.* our targeted estimate of \mathbf{D}) is either *zero* or *close to zero* (probably because the true \mathbf{D} is so). A traditional way to anticipate such a situation usually consists in replacing (S:2.2a) with a mix of both relative and absolute error controls:

$$\|\hat{\mathbf{D}}^{(t+1)} - \hat{\mathbf{D}}^{(t)}\|^{(M)} \leq \delta_1 \cdot \|\hat{\mathbf{D}}^{(t+1)}\|^{(M)} \quad \text{or} \quad \|\hat{\mathbf{D}}^{(t+1)} - \hat{\mathbf{D}}^{(t)}\|^{(M)} \leq \delta_{11}, \quad (\text{S:2.3})$$

where δ_{11} is an absolute error tolerance upper bound.

But because it is generally difficult to set an appropriate value for δ_{11} in a given context, here we instead use the following context dependent alternative to the second inequality in (S:2.3):

$$\|\hat{\mathbf{D}}^{(t+1)} - \hat{\mathbf{D}}^{(t)}\|^{(M)} \leq 0.5 \cdot \|\hat{\mathbf{D}}^{(t+1)}\|^{(M)} \quad \text{and} \quad \|\mathbf{Z}\hat{\mathbf{G}}^{(t+1)}\mathbf{Z}^T\|^{(M)} \leq 0.01 \cdot \hat{\sigma}_\varepsilon^{2(t+1)}, \quad (\text{S:2.4})$$

the latter inequality aiming at detecting a high probability of zero cluster random effects in the given data set. Hence, the final stopping criterion used in our 3S algorithms is rather:

$$((\text{S:2.2a}) \text{ or } (\text{S:2.4})) \text{ and } (\text{S:2.2b}) \text{ and } (\text{S:2.2c}). \quad (\text{S:2.5})$$

The rightmost inequality in (S:2.4) tries to test whether the contribution to the response covariance matrix \mathbf{V} of the random effects covariance matrix \mathbf{G} is (at least) *two orders of magnitude smaller* than that of the residual errors variance σ_ε^2 . This occurring gives an indication that these random effects either are not really there or are insignificant, which is slowing down the convergence of the $\hat{\mathbf{D}}^{(t)}$ iterates, and thus the algorithm. So when (S:2.4), (S:2.2b) and (S:2.2c) are satisfied, but not (S:2.2a), we stop with `code.Stop = OK:2`, inviting the user to view this as a *partial success* scenario.

5) `code.Warning = 3: RE might not be there, anyway...`

A 3S algorithm being outright successful (*i.e.* termination with `code.Stop = OK:1`) does not imply that there are, indeed, random effects in the provided data. So it is useful to try to detect their presence or absence anyway, convergence of the algorithm or not. Our way of doing so is to always test the rightmost inequality in (S:2.4) after exiting a 3S algorithm. If satisfied, a warning is issued to the user: *despite the observed convergence (if any), with the given clustering variable(s), random effects might actually not be there, or seem insignificant*. An ad-

vantage in that way of trying to detect the zero random effects scenario is that it works whatever the dimension r of these ones in each cluster.

In any event, we always include, in the output from the algorithm, the value of the last ratio:

$$\hat{\rho} = \frac{\|\mathbf{Z}\hat{\mathbf{G}}^{(t)}\mathbf{Z}^T\|^{(M)}}{\hat{\sigma}_\varepsilon^{2(t)}}, \quad (\text{S:2.6a})$$

estimating

$$\rho = \frac{\|\mathbf{Z}\mathbf{G}\mathbf{Z}^T\|^{(M)}}{\sigma_\varepsilon^2}. \quad (\text{S:2.6b})$$

The latter can be called the *random effects ratio* in the considered LMM for the given data set, and can be roughly interpreted as the coefficient of random effects really present in the data:

- 1) the higher above the number 1 the ratio ρ is, the more effective such a presence is;
- 2) the smaller below 1 the ratio ρ is, the more doubt can be cast about that presence.

But, for a 2-level LMM (2), note that the block diagonal structure of \mathbf{Z} implies:

$$\|\mathbf{Z}\mathbf{G}\mathbf{Z}^T\|^{(M)} = \max_{1 \leq j \leq m} \|\mathbf{Z}_j \mathbf{D} \mathbf{Z}_j^T\|^{(M)}, \quad \|\mathbf{Z}\hat{\mathbf{G}}^{(t)}\mathbf{Z}^T\|^{(M)} = \max_{1 \leq j \leq m} \|\mathbf{Z}_j \mathbf{D}^{(t)} \mathbf{Z}_j^T\|^{(M)}. \quad (\text{S:2.7})$$

Remark S2.2. *One may wonder: why not simply use a statistical test to decide whether there are cluster random effects or not in the given data? The answer is twofold. First, theoretically valid statistical tests are hard to develop for LMMs, more so in our context where we aim at fitting these models without imposing any parametric distributional assumptions on the data. Secondly, as seen above, the situation of no cluster random effects actually first impacts negatively on the convergence of a 3S algorithm (as well as other iterative algorithms for LMM fitting), often severely slowing it down or even flat out stalling it. Hence in that scenario, it is difficult to even compute reliable estimates on which to perform statistical tests. Here, we wanted to develop practical numerical tests aimed at detecting it during the run of the algorithm. And such a detection is of the utmost practical importance since, if well done, it may allow to decide whether LMM modeling even makes sense or not for a given data set. Nevertheless, upon exiting a 3S algorithm, we can perform a bootstrap test to decide whether the ratio ρ is, or not, significantly greater than zero for the given data set.*

S2.3.3. code.Warning = 4: Almost Constant Sequence of $\hat{\beta}^{(t)}$ Iterates

For many data sets, we observed that the sequence of $\hat{\beta}^{(t)}$ iterates produced by any of our 3S algorithms is nearly constant. This has to do with our way of initializing that sequence by the OLS in an LM fit of the data, i.e. an LMM with no RE. With that initialization, if the sequence ($\hat{\beta}^{(t)}$) exhibits a constant trend from

the outset, one may suspect that the LM is probably better suited for the data set at hand than a pure LMM, *i.e.* an LMM with nonzero RE. And this is likely so because either there are no cluster random effects or these are nonsignificant or would require a bigger sample size to be detected. Hence, whatever the stoppage condition (*convergence or not*) of the iterative algorithm, if for all iterations carried, the $\hat{\beta}^{(t)}$ iterates were all only negligibly different from their OLS starting point $\hat{\beta}^{(0)}$, we conclude that we might be in that *zero RE* scenario, meaning that *what was supposed to be cluster random effects actually behave, more or less, the same from one cluster to the other.*

More specifically, we always perform the following test after termination of any of our 3S algorithms (successfully or not):

$$\text{for } t = 1(1)T, \quad \|\hat{\beta}^{(t)} - \hat{\beta}^{(0)}\| \leq \delta_2 \|\hat{\beta}^{(0)}\|, \quad (\text{S:2.8})$$

where $t = T$ is the iteration at which the algorithm was stopped, and δ_2 is the same as in (S:2.2b). If (S:2.8) is satisfied, then the user is informed that: *may be the LM is better suited for her/his data set.* The recommendation then is either to use the LM, or redesign the clustering variable because it might have been poorly chosen in the first place.

S2.3.4. code.Stop = KO:5: Authorized Maximum Number of Iterations Reached without Convergence

As is customary for iterative algorithms expected to stop at satisfaction of a convergence criterion, we must also anticipate the possibility that, for the provided data set, a 3S algorithm either does not actually converge or does so slowly. For that, we set, as is usual, a maximum number max.Its (default = 200) of authorized iterations in the algorithm. If that number is reached without the convergence stopping criterion satisfied, the algorithm is terminated with code.Stop = 5, and we inform the user why.

S2.3.5. Complement: Authorized Minimum Number of Iterations

Somewhat untraditionally, but to limit the risk of an optimistic premature successful termination, we impose that a 3S algorithm always performs a minimum number min.Its (default = 10) before starting to check for a success in convergence.

S3. Two Lemmas Useful for Proving Theorem 7

We admit the following two lemmas (the second one can be proved using multivariable differential tools presented in Graham (1981)):

Lemma S3.1. *Let $S \in \mathcal{M}_n(\mathbb{R})$. Then, one has: $2S - \text{diag}(S) = 0 \Leftrightarrow S = 0$.*

Lemma S3.2. *Let $x \in \mathbb{R}^n$ and $N \in \mathcal{M}_n(\mathbb{R})$, SPD, such that*

$$N = ZQZ^T + \alpha I_n, \quad \text{with } Q \in \mathcal{M}_r(\mathbb{R}), \text{ symmetric, } \alpha \in \mathbb{R} \text{ and } Z \in \mathcal{M}_{n,r}(\mathbb{R}).$$

If the upper triangular part of Q has functionally independent elements and which are independent from α , whereas x and Z are independent from Q and α , then:

$$\frac{\partial(\mathbf{x}^T N^{-1} \mathbf{x})}{\partial Q} = \text{diag}(W) - 2W, \tag{S:3.1a}$$

$$\frac{\partial \log |N|}{\partial Q} = 2B - \text{diag}(B), \tag{S:3.1b}$$

$$\frac{\partial(\mathbf{x}^T N^{-1} \mathbf{x})}{\partial \alpha} = -\|N^{-1} \mathbf{x}\|^2, \tag{S:3.1c}$$

$$\frac{\partial \log |N|}{\partial \alpha} = \text{tr}(N^{-1}), \tag{S:3.1d}$$

where $W = Z^T N^{-1} \mathbf{x} \mathbf{x}^T N^{-1} Z$, $B = Z^T N^{-1} Z$, and if A is a square matrix, then $\text{diag}(A)$ is the diagonal matrix with same main diagonal as A .

S4. 3S Fitting Algorithms for 2-Level LMMs with u.ho. Errors and Diagonal Covariance Matrix for the Cluster Random Effects, and ANOVA LMMs

S4.1. 3S Algorithms for 2-Level LMMs with u.ho. Errors Assuming a Diagonal Covariance Matrix for the Cluster Random Effects

We present, here, the equivalence of Algorithms 3S-A1-V1 and 3S-A1-V2 respectively for 2-level LMMs with u.ho. errors, assuming a diagonal covariance matrix for the cluster random effects vector, that is:

$$D = \text{diag}(\sigma_1^2, \dots, \sigma_r^2).$$

They carry respectively the code names 3S-A1-V1-d and 3S-A1-V2-d.

Algorithm 3S-A1-V1-d. Estimation of β , σ_ε^2 , $D = \text{diag}(\sigma_1^2, \dots, \sigma_r^2)$ and prediction of U_1, \dots, U_m in a 2-level LMM with u.ho. errors. Version 1

1) **Initialization:** At iteration 0, we estimate β , σ_ε^2 , D and predict U_1, \dots, U_m , as follows:

a) As in Algorithm 3S-A1-V1;

b) $\hat{\sigma}_k^{2(0)} = \frac{1}{m} \sum_{j=1}^m (\hat{u}_{j,k}^{(0)})^2$, $k = 1, \dots, r$, with $\hat{u}_{j,k}^{(0)}$ the k^{th} element of the vector $\hat{\mathbf{u}}_j^{(0)}$;

c) $\hat{D}^{(0)} = \text{diag}(\hat{\sigma}_1^{2(0)}, \dots, \hat{\sigma}_r^{2(0)})$;

2) **The iterative process:** Given $\hat{\beta}^{(t)}$, $\hat{\sigma}_\varepsilon^{2(t)}$, $\hat{\mathbf{u}}_1^{(t)}, \dots, \hat{\mathbf{u}}_m^{(t)}$, hence $\hat{\mathbf{u}}^{(t)} = (\hat{\mathbf{u}}_1^{(t)T}, \dots, \hat{\mathbf{u}}_m^{(t)T})^T$, and $\hat{D}^{(t)} = \text{diag}(\hat{\sigma}_1^{2(t)}, \dots, \hat{\sigma}_r^{2(t)})$ from iteration t , we obtain estimates and predictions at iteration $t+1$ as follows:

a) As in Algorithm 3S-A1-V1;

b) $\hat{\sigma}_k^{2(t+1)} = \frac{1}{m} \sum_{j=1}^m \left[(\hat{\mathbf{V}}_j^{*(t+1)})_{kk} + (\hat{u}_{j,k}^{(t+1)})^2 \right]$, $k = 1, \dots, r$;

c) $\hat{D}^{(t+1)} = \text{diag}(\hat{\sigma}_1^{2(t+1)}, \dots, \hat{\sigma}_r^{2(t+1)})$;

3) **Stopping criterion:** We assume convergence when all the 3 inequalities

$$\|\hat{\beta}^{(t+1)} - \hat{\beta}^{(t)}\| \leq \delta_1 \cdot \|\hat{\beta}^{(t+1)}\|, \tag{S:4.1a}$$

$$\|\hat{\beta}^{(t+1)} - \hat{\beta}^{(t)}\| \leq \delta_2 \cdot \|\hat{\beta}^{(t+1)}\|, \tag{S:4.1b}$$

$$\left| \hat{\sigma}_\varepsilon^{2(t+1)} - \hat{\sigma}_\varepsilon^{2(t)} \right| \leq \delta_3 \cdot \hat{\sigma}_\varepsilon^{2(t+1)}, \quad (\text{S:4.1c})$$

are satisfied, where $\delta_1, \delta_2, \delta_3$ are relative tolerance levels set in $(0, 1)$, and $\hat{\sigma} = (\hat{\sigma}_1^2, \dots, \hat{\sigma}_r^2)^\top$. Otherwise, repeat Step 2 with $t \leftarrow t+1$.

4) **Extracting estimates.** At convergence, take $\hat{\beta} = \hat{\beta}^{(t+1)}$, $\hat{\sigma}_\varepsilon^2 = \hat{\sigma}_\varepsilon^{2(t+1)}$, $\hat{D} = \hat{D}^{(t+1)}$ as estimates of β , σ_ε^2 , D . Also, take $\hat{u}_1^{(t+1)}, \dots, \hat{u}_m^{(t+1)}$ as predictions of U_1, \dots, U_m .

Algorithm 3S-A1-V2-d. Estimation of β , σ_ε^2 , $D = \text{diag}(\sigma_1^2, \dots, \sigma_r^2)$ and prediction of U_1, \dots, U_m in a 2-level LMM with u.h.o. errors. Version 2

- 1) **Initialization:** The same as in Algorithm 3S-A1-V1-d;
- 2) **The iterative process.** Given $\hat{\beta}^{(t)}$, $\hat{\sigma}_\varepsilon^{2(t)}$, $\hat{u}_1^{(t)}, \dots, \hat{u}_m^{(t)}$, hence $\hat{u}^{(t)} = (\hat{u}_1^{(t)\top}, \dots, \hat{u}_m^{(t)\top})^\top$, and $\hat{D}^{(t)} = \text{diag}(\sigma_1^{2(t)}, \dots, \sigma_r^{2(t)})$ from iteration t , we obtain estimates and predictions at iteration $t+1$ as follows:
 - a) As in Algorithm 3S-A1-V2;
 - b) As in Algorithm 3S-A1-V1-d;
 - 3) **Stopping criterion:** As in Algorithm 3S-A1-V1-d. Otherwise, repeat Step 2 with $t \leftarrow t+1$.
 - 4) **Extracting estimates.** As in Algorithm 3S-A1-V1-d.

S4.2. 3S Algorithms for ANOVA LMMs

We present the adaptation of Algorithms 3S-A1-V1 and 3S-A1-V2 respectively for ANOVA LMMs. That is, the case where

$$D = G = \text{diag}(\sigma_1^2 \mathbf{I}_{m_1}, \dots, \sigma_g^2 \mathbf{I}_{m_g}),$$

where, for $k = 1, \dots, g$, each variance σ_k^2 corresponds to the k th categorical variable with m_k categories, with g the total number of categorical variables with random effects. They carry respectively the code names 3S-ANOVA-A1-V1 and 3S-ANOVA-A1-V2.

Algorithm 3S-ANOVA-A1-V1. Estimation of β , σ_ε^2 , $D = G$ and prediction of U in a 2-level ANOVA LMM: Version 1

- 1) **Initialization:** At iteration 0, we estimate β , σ_ε^2 , G and predict U , as follows:
 - a) As in Algorithm 3S-A1-V1;
 - b) $\hat{u}^{(0)} = (\mathbf{Z}^\top \mathbf{Z})^{-1} \mathbf{Z}^\top (\mathbf{y} - \mathbf{X} \hat{\beta}^{(0)})$, the OLS estimate of U in the linear model $\mathbf{Y}^{*(0)} = \mathbf{Z} \mathbf{U} + \varepsilon$, with U considered as fixed parameter in the model, \mathbf{Z} the design matrix, and $\mathbf{Y}^{*(0)} = \mathbf{Y} - \mathbf{X} \hat{\beta}^{(0)}$ as the response vector;
 - c) $\hat{\sigma}_\varepsilon^{2(0)} = \frac{\|\mathbf{y} - \mathbf{X} \hat{\beta}^{(0)} - \mathbf{Z} \hat{u}^{(0)}\|^2}{(n-p)}$;
 - d) $\hat{\sigma}_k^{2(0)} = \frac{1}{m_k} \|\hat{u}_k^{(0)}\|^2 = \frac{1}{m_k} \sum_{l=1}^{m_k} (\hat{u}_{k,l}^{(0)})^2$, for $k = 1, \dots, g$, with $\hat{u}^{(0)} = (\hat{u}_1^{(0)\top}, \dots, \hat{u}_g^{(0)\top})^\top$ split according to the g categorical variables that form the boolean matrix \mathbf{Z} , such that $\hat{u}_k^{(0)} = (\hat{u}_{k,1}^{(0)}, \dots, \hat{u}_{k,m_k}^{(0)})^\top$ is the m_k -vector of random effects of the k th categorical variable.
 - e) $\hat{G}^{(0)} = \text{diag}(\hat{\sigma}_1^{2(0)} \mathbf{I}_{m_1}, \dots, \hat{\sigma}_g^{2(0)} \mathbf{I}_{m_g})$;

2) **The iterative process.** Given $\hat{\beta}^{(t)}$, $\hat{\sigma}_\varepsilon^{2(t)}$, $\hat{u}^{(t)}$, and $\hat{G}^{(t)} = \text{diag}(\hat{\sigma}_1^{2(t)} I_{m_1}, \dots, \hat{\sigma}_g^{2(t)} I_{m_g})$ from iteration t , we obtain estimates and predictions at iteration $t+1$ as follows:

a) As in Algorithm 3S-A1-V1;

b) $\hat{V}^{*(t+1)} = (\hat{\sigma}_\varepsilon^{-2(t+1)} Z^T Z + \hat{G}^{(t-1)})^{-1}$;

c) $\hat{\sigma}_k^{2(t+1)} = \frac{1}{m_k} \sum_{l=1}^{m_k} \left[(\hat{V}_{kk}^{*(t+1)})_{ll} + (\hat{u}_{k,l}^{(t+1)})^2 \right]$, $k = 1, \dots, g$, where $\hat{V}_{kk}^{*(t+1)}$ is the k^{th} diagonal block of $\hat{V}^{*(t+1)}$ corresponding to the m_k categories of the k^{th} categorical random effects variable;

d) $\hat{G}^{(t+1)} = \text{diag}(\hat{\sigma}_1^{2(t+1)} I_{m_1}, \dots, \hat{\sigma}_g^{2(t+1)} I_{m_g})$;

3) **Stopping criterion:** We assume convergence when all the 3 inequalities

$$\|\hat{\beta}^{(t+1)} - \hat{\beta}^{(t)}\| \leq \delta_1 \cdot \|\hat{\beta}^{(t+1)}\|, \tag{S:4.2a}$$

$$\|\hat{\beta}^{(t+1)} - \hat{\beta}^{(t)}\| \leq \delta_2 \cdot \|\hat{\beta}^{(t+1)}\|, \tag{S:4.2b}$$

$$|\hat{\sigma}_\varepsilon^{2(t+1)} - \hat{\sigma}_\varepsilon^{2(t)}| \leq \delta_3 \cdot \hat{\sigma}_\varepsilon^{2(t+1)}, \tag{S:4.2c}$$

are satisfied, where $\delta_1, \delta_2, \delta_3$ are relative tolerance levels set in $(0,1)$, and $\hat{\beta} = (\hat{\sigma}_1^2, \dots, \hat{\sigma}_g^2)^T$. Otherwise, repeat Step 2 with $t \leftarrow t+1$.

4) **Extracting estimates.** At convergence, take $\hat{\beta} = \hat{\beta}^{(t+1)}$, $\hat{\sigma}_\varepsilon^2 = \hat{\sigma}_\varepsilon^{2(t+1)}$, $\hat{G} = \hat{G}^{(t+1)}$ as estimates of β , σ_ε^2 , G . Also, take $\hat{u}^{(t+1)}$ as predictions of U .

Algorithm 3S-ANOVA-A1-V2. Estimation of β , σ_ε^2 , $D = G$ and prediction of U in a 2-level ANOVA LMM: Version 2

1) **Initialization:** The same as in Algorithm 3S-ANOVA-A1-V1;

2) **The iterative process.** Given $\hat{\beta}^{(t)}$, $\hat{\sigma}_\varepsilon^{2(t)}$, $\hat{u}^{(t)}$, and $\hat{G}^{(t)} = \text{diag}(\hat{\sigma}_1^{2(t)} I_{m_1}, \dots, \hat{\sigma}_g^{2(t)} I_{m_g})$ from iteration t , we obtain estimates and predictions at iteration $t+1$ as follows:

a) As in Algorithm 3S-A1-V2, but removing the precision about $\hat{G}^{(t)}$;

b) $A^{(t+1)} = (Z^T Z + \hat{\sigma}_\varepsilon^{2(t+1)} \hat{G}^{(t-1)})^{-1}$;

c) $\hat{u}^{(t+1)} = A^{(t+1)} Z^T (y - X \hat{\beta}^{(t+1)})$;

d) $\hat{V}^{*(t+1)} = \hat{\sigma}_\varepsilon^{2(t+1)} A^{(t+1)}$;

e) As in Algorithm 3S-ANOVA-A1-V1;

3) **Stopping criterion:** As in Algorithm 3S-ANOVA-A1-V1. Otherwise, repeat Step 2 with $t \leftarrow t+1$.

4) **Extracting estimates.** As in 3S-ANOVA-A1-V1.

55. Parameters of the Distributions Used for the Simulations

In this section, we detail the distributions and parameters used for the simulations that yield the results in Section 6.1. As a general rule, the values of most of these parameters were themselves generated randomly before launching the simulations.

55.1. Simulation of Sample Items, Clusters, Fixed and Random Effects Covariates

We assume that we have a population comprising $m = 10$ clusters in the same

proportions of 1/10. To get a sample item, we first sampled its cluster, then its fixed and random effects covariates, respectively in \mathbb{R}^4 and \mathbb{R}^3 , assuming that in each cluster j , the vector X of fixed effects covariates followed $\mathcal{N}_4(\mathbf{m}_j^X, \mathbf{I}_4)$. Likewise, in each cluster j , the vector Z of random effects covariates followed $\mathcal{N}_3(\mathbf{m}_j^Z, 0.7\mathbf{I}_3)$. We used:

- $\mathbf{m}_1^X = (-2, 1, -2, 2)^\top$, $\mathbf{m}_2^X = (-4, -4, -4, 1)^\top$, $\mathbf{m}_3^X = (4, -1, 3, 2)^\top$,
 $\mathbf{m}_4^X = (-1, 3, 0, -2)^\top$, $\mathbf{m}_5^X = (-3, 3, 2, -2)^\top$, $\mathbf{m}_6^X = (2, -1, 3, -3)^\top$,
 $\mathbf{m}_7^X = (-3, -3, 4, 4)^\top$, $\mathbf{m}_8^X = (-4, 1, 4, -3)^\top$, $\mathbf{m}_9^X = (-1, 2, 0, -3)^\top$,
 $\mathbf{m}_{10}^X = (-1, 0, 0, 4)^\top$;
- $\mathbf{m}_1^Z = (-2, 0, -3)^\top$, $\mathbf{m}_2^Z = (3, 1, 1)^\top$, $\mathbf{m}_3^Z = (0, -1, -1)^\top$, $\mathbf{m}_4^Z = (2, 0, -2)^\top$,
 $\mathbf{m}_5^Z = (-3, 2, 3)^\top$, $\mathbf{m}_6^Z = (3, -3, 0)^\top$, $\mathbf{m}_7^Z = (-2, 3, 2)^\top$, $\mathbf{m}_8^Z = (2, 3, -3)^\top$,
 $\mathbf{m}_9^Z = (0, 3, 1)^\top$, $\mathbf{m}_{10}^Z = (-2, 2, 0)^\top$.

In addition, no intercept was included, be it in the fixed part or the random part of the simulated LMM, while we took $\beta = (-1.5, 0.7, 0, 2.3)^\top \in \mathbb{R}^4$ as vector of fixed effects parameters.

S5.2. Parameters of the Distributions for Simulating Random Effects and Residual Errors

To simulate the cluster random effects in \mathbb{R}^3 , we examined 5 distributions, while for the residual errors, we examined 3, hence, testing in total, $5 \times 3 = 15$ scenarios.

S5.2.1. Distributions Tested for the Random Effects

Case 1: Non-degenerate Gaussian distribution in \mathbb{R}^3 with mean vector and covariance matrix:

$$\mathbf{m} = \mathbf{0}, \quad \Sigma = \begin{pmatrix} 32.784612 & 28.546477 & -9.726209 \\ 28.546477 & 44.499170 & -3.915344 \\ -9.726209 & -3.915344 & 27.435629 \end{pmatrix}.$$

Case 2: A degenerate Gaussian distribution in \mathbb{R}^3 with mean vector and a rank one covariance matrix:

$$\mathbf{m} = \mathbf{0}, \quad \Sigma = \begin{pmatrix} 2.5820307 & -11.232197 & -0.46004207 \\ -11.2321974 & 48.861640 & 2.00124781 \\ -0.4600421 & 2.001248 & 0.08196599 \end{pmatrix}.$$

Case 3: As in Case 2, but with mean vector and a rank two covariance matrix:

$$\mathbf{m} = \mathbf{0}, \quad \Sigma = \begin{pmatrix} 10.19483 & 6.661080 & -17.54595 \\ 6.66108 & 5.462537 & -10.91097 \\ -17.54595 & -10.910967 & 30.47329 \end{pmatrix}.$$

Case 4: Mixture of 3 Gaussian distributions in \mathbb{R}^3 with vector of mixing proportions $\pi = (1/6, 1/3, 1/2)^\top$, common covariance matrix $\Sigma = 2\mathbf{I}_3$ and respective vector means:

$$\mathbf{m}_1 = (0.2815777, 1.8254139, -0.5799608)^\top,$$

$$\mathbf{m}_2 = (-0.16064371, 0.01403464, -1.63159078)^\top,$$

$$\mathbf{m}_3 = (0.01323658, -0.61782773, 1.28104746)^T.$$

Case 5: A discrete distribution in \mathbb{R}^3 with vector of respective probabilities $\pi = (1/6, 1/3, 1/2)^T$ on the 3 mass points:

$$\mathbf{u}_1 = \begin{pmatrix} 0.7043073 \\ 0.7533598 \\ -0.7370090 \end{pmatrix}, \quad \mathbf{u}_2 = \begin{pmatrix} 0.6410311 \\ 0.8449623 \\ -0.7769852 \end{pmatrix}, \quad \mathbf{u}_3 = \begin{pmatrix} -3.3851412 \\ 0.3501683 \\ 0.8949348 \end{pmatrix}.$$

S5.2.2. Distributions Tested for the Residual Errors

Case 1: Univariate Gaussian distribution $\mathcal{N}(0, 2.25)$.

Case 2: Mixture of 4 Gaussian distributions in \mathbb{R} of variance $\sigma^2 = 0.9613748$ and respective means:

$$\mathbf{m}_1 = 1.293309, \quad \mathbf{m}_2 = -1.912166, \quad \mathbf{m}_3 = 0.4601711, \quad \mathbf{m}_4 = 0.690338,$$

and vector of mixing proportions $\pi = (1/8, 1/4, 1/2, 1/8)^T$.

Case 3: A discrete distribution in \mathbb{R} with the 4 mass points:

$$u_1 = 1.532882, \quad u_2 = 1.332857, \quad u_3 = -1.494380, \quad u_4 = 1.778924$$

of respective probabilities $\pi = (1/8, 1/4, 1/2, 1/8)^T$.

S6. More Analysis of the Blackmore Data

One can also consider fitting an LMM to the Blackmore data where, in addition to having a random slope for age as in (44), we also have a random slope for the interaction age8*group, thus the following LMM:

$$\begin{aligned} \log_2 \text{exercise} \sim & \text{fe}(1 + \text{age8} + \text{group} + \text{age8} * \text{group}) \\ & + \text{re}(1 + \text{age8} + \text{age8} * \text{group}) + \text{Gr}(\text{subject}) \end{aligned} \tag{S:6.1}$$

The results for this model using Algorithm 3S-A1-V1 are presented in **Table S1**. In addition to the parameters which are already significant in the LMM (44), the estimate of the variance $\sigma_{\text{age8*group}}^2$ of the random effect age8*group in the LMM (S:6.1) is also significantly greater than zero, confirming that a random effect on age8*group is also probable. This suggests that there is a likely variation of the interaction between age and group across the girls population in this study.

Table S1. Results of algorithm 3S-A1-V1 fitting the LMM (S:6.1) to the Blackmore data.

parameters	est.	[2.5% – 97.5%]	bias	std.dev	$\sqrt{\text{MSE}}$	t-stat	p-value
β_{int}	-0.2759	[-0.6373, 0.09259]	0.0008132	0.1852	0.1852	-1.490	1.363e - 01
β_{age8}	0.06383	[0.008403, 0.1242]	0.001310	0.03032	0.03035	2.103	3.547e - 02
$\beta_{\text{gr:pat}}$	-0.3546	[-0.8601, 0.1041]	-0.0005366	0.2402	0.2402	-1.476	1.399e - 01
$\beta_{\text{age8*gr:pat}}$	0.2402	[0.1585, 0.3151]	-0.002138	0.03931	0.03936	6.103	1.041e - 09
σ_{int}^2	2.094	[1.693, 2.652]	0.07693	0.2418	0.2537	8.254	1.532e - 16
σ_{age8}^2	0.02712	[0.01910, 0.04529]	0.003153	0.006638	0.007349	3.690	2.239e - 04
$\sigma_{\text{age8*gr:pat}}^2$	0.01366	[0.01123, 0.01957]	0.001324	0.002206	0.002573	5.309	1.101e - 07
$\text{CORR}_{\text{int, age8}}$	-0.1995	[-0.4627, -0.04233]	-0.06773	0.3306	0.3586	-1.551	1.209e - 01
$\text{CORR}_{\text{int, age8*gr:pat}}$	-0.2141	[-0.4588, 0.01548]	-0.01022	0.3530	0.3536	-1.713	8.677e - 02
$\text{CORR}_{\text{age8, age8*gr:pat}}$	-0.2666	[-0.3945, 0.1250]	0.1443	0.3628	0.4419	-1.365	1.722e - 01
σ_{ϵ}^2	1.530	[1.2860, 1.662]	-0.05381	0.09631	0.1103	13.870	1.017e - 43
ρ	2.497	[1.841, 3.986]	0.3309	0.5463	0.6388	3.909	9.252e - 05

int, gr, pat: respective shortcuts for intercept, group, patient.

In contrast, we signal that Gaussian ML fitting of the LMM (S:6.1) to the Blackmore data failed to converge with the default settings in the lmer function in the lme4 R package.

References

- Anderson, E. and ten others (1999). *LAPACK Users' Guide, 3rd Edition*, SIAM, Philadelphia, Available on-line at http://www.netlib.org/lapack/lug/lapack_lug.ht.
- Demmel, J. W. (1997). *Applied Numerical Linear Algebra*, SIAM, Philadelphia.
- Dongarra, J. J., Bunch, J. R., Moler, C. B. and Stewart, G. W. (1978). *LINPACK Users Guide*, SIAM, Philadelphia, PA, USA.
- Graham, A. (1981). *Kronecker Products and Matrix Calculus: with Applications*, Ellis Horwood, Chichester, West Sussex, England.