# Mathematics

*Research article*

# A new test for detecting specification errors in Gaussian linear mixed-effects models

**Jairo A. Angel[1,*], Francisco M.M. Rocha[2], Jorge I. Vélez[3] and Julio M. Singer[4]**

[1] Department of Mathematics and Statistics, University of Córdoba, Montería, Colombia
[2] Paulista School of Economics and Business Politics, Federal University of São Paulo, Brazil
[3] Department of Industrial Engineering, Universidad del Norte, Barranquilla, Colombia
[4] Department of Statistics, Instituto de Matemática e Estatística, University of São Paulo, São Paulo, Brazil

* **Correspondence:** E-mail: jaangelguzman@correo.unicordoba.edu.co.

**Abstract:** Linear mixed-effects models (LMEMs) are widely used in medical, engineering, and social applications. The accurate specification of the covariance matrix structure within the error term is known to impact the estimation and inference procedures. Thus, it is crucial to detect the source of errors in LMEMs specifications. In this study, we propose combining a user-friendly computational test with an analytical method to visualize the source of errors. Through statistical simulations under different scenarios, we evaluate the performance of the proposed test in terms of the Power and Type I error rate. Our findings indicate that as the sample size $n$ increases, the proposed test effectively detects misspecification in the systematic component, the number of random effects, the within-subject covariance structure, and the covariance structure of the error term in the LMEM with high Power while maintaining the nominal Type I error rate. Finally, we show the practical usefulness of our proposed test with a real-world application.

**Keywords:** linear mixed-effects model; misspecification; Type I error; simulation
**Mathematics Subject Classification:** 62F05, 62J05, 62J20

## 1. Introduction

In linear mixed-effects models (LMEMs), the assumptions that concern the structure of the covariance matrix of the response variable are presumed to be adequately specified. This involves specifying the mean structure, the covariance matrix structure, and the distribution pattern of the covariance matrix. These elements define both the covariance among the individuals and those associated with the vector of random effects. However, verifying these assumptions in LMEMs can be

challenging due to the complexity of the data structure and the presence of two sources of error within the model.

Diagnostic tests are crucial for detecting misspecifications errors in LMEMs, as they identify potential model misrepresentations through analytical methods. Some authors have proposed tests to detect misspecification when formulating LMEMs. For instance, [3] used a strategy to evaluate the random effect distribution via a parametric Bootstrap for small samples in the case of mixed models, generalized linear models, and non-linear mixed models. For this purpose, they used an asymptotic test based on a gradient function

Through statistical simulations, [10] showed that the increase in Type II error was consistent with the effect of specification error on the distribution of the random effect for generalized LMEMs. On the other hand, [7] developed a diagnostic method by performing data reconstruction, to detect misspecification for generalized LMEMs. The author proposed a theoretical justification of the method and investigated the behavior of this method via a simulation in finite samples. In these simulations, the author compared a model without a specification error with a model with a misspecification in its fixed part.

Several techniques have emerged to discern errors in either the distribution of the error term or the distribution of random effects within statistical models. These components have been assessed by different researchers to ensure model integrity and draw accurate conclusions [5,9,16]. Specifically, [8] proposed a test to analyze specification errors in the distribution of the random error term and the random effects, while [7] presented a test that allowed them to identify errors in the specification of the distribution of the random effects.

For LMEMs, exploratory techniques have also been proposed to identify the sources of errors. For instance, [18] proposed a set of graphical and analytical techniques, based on three types of residuals (i.e., marginal, conditional, and random effects) for the diagnosis of the intra-unit sample covariance matrix in repeated measure studies, as well as graphical tools to analyze violations of the error structure in LMEMs. To identify the number of random effects in an LMEM, the recommended exploratory methods based on individual and mean profile analyses [15].

Despite their utility, the tests previously described fall short in pinpointing the origins or specific sources of these errors. Hence, int his paper we propose combining a user-friendly computational test with an analytical method to visualise the source of errors, while considering the methodology suggested by [15]. However, our focus does not compare our approach with other approaches in the literature to determine the superior power behaviour. Instead, we concentrate on integrating two approaches: a formal test and a graphical diagnostic tool. This combination allows for the detection of model misspecifications and the identification of source of errors within the model. Hence, the test detects when the model is misspecified and the analytical method allows us to visualise the source of the misspecification.

This article is organized as follows: in Section 2, the Gaussian LMEM is outlined; in Section 3, the proposed test is described; in Section 4, the simulation study and the different scenarios are described in detail, and in Section 5 the results are reported; subsequently, in Section 6, we illustrate the usefulness of the proposed test using real data, and present the diagnostic graphical tools to visualize the source of the misspecification; and finally, the conclusions are recommendations presented, and further areas of research are discussed.

## 2. The Gaussian linear-mixed effects model

A general form of the Gaussian LMEM is

$$y_i = X_i\beta + Z_i b_i + e_i, \qquad i = 1, \ldots n, \tag{2.1}$$

where $y_i = (y_{i_1}, y_{i_2}, ..., y_{i_{m_i}})^\top$ represents the vector of the $m_i$ observations recorded for the $i$-th sample unit, $\beta = (\beta_1, \ldots, \beta_p)^\top$ denotes the vector of either the location parameters or the fixed effects, $X_i$ is the matrix corresponding to the specification of the fixed terms, $b_i = (b_{1i}, \ldots, b_{qi})^\top$ represents the vector of random effects, $Z_i$ is the matrix corresponding to the specification of the vector of random effects and $e_i$ represents the vector of random errors. By definition, $b_i \sim \mathcal{N}_{m_i}[0, G(\theta)]$ and $e_i \sim \mathcal{N}_{m_i}[0, R_i(\theta)]$ are assumed to be independent.

The vector $\theta = (\theta_1, \ldots \theta_k)^\top$ contains all non-redundant components (parameters) of the covariance matrix of vector $b_i$. The vector $\phi = (\beta^\top, \theta^\top)^\top (s \times 1)$ represents the vector of all the parameters in (2.1). It follows that the covariance matrix of the vector $y_i$ can be written as

$$\mathbb{V}ar(y_i) = \mathbb{V}_i(\theta) = \mathbb{V}_i = Z_i G Z_i^\top + R_i, \tag{2.2}$$

where $G = G(\theta)$ and $R_i = R_i(\theta)$. In summary,

$$y_i \sim \mathcal{N}_{m_i} [X_i\beta, \mathbb{V}_i]. \tag{2.3}$$

The most commonly used estimation methods for the parameters in model (2.1) are the Maximum Likelihood Estimation (MLE) method and the Restricted MLE (RMLE). For additional details on the MLE and the RMLE, see [12] and [14], respectively.

The maximum likelihood methodology yields unbiased estimators for the fixed effects, though it introduces bias in the estimators for the random effects. This bias stems from disregarding the loss of the degrees of freedom during the estimation of the fixed terms. Consequently, it also results in biased estimators for the parameters of the intra-unit covariance matrix.

The random vectors $y_1, \ldots, y_n$ are independent with a distribution given by (2.3). The probability density function associated to each vector $y_i$ is denoted as $f(y_i; \phi)$. Considering the vector

$$y = \left(y_i^T, ..., y_n^T\right)^T,$$

and considering the probability density function of the random variables $y_i$, the likelihood function of $\phi$ is as follows:

$$L(\phi; y) = \prod_{1 \le i \le n} f(y_i; \phi) = \prod_{1 \le i \le n} \int_{\mathcal{R}^{q_i}} f(y_i; \phi \,|\, b_i) f(b_i; \phi) db_i. \tag{2.4}$$

where $\mathcal{R}^{q_i}$ is the $q_i$-dimensional space of the vector $b_i$. It follows that the logarithm of (2.4) is represented by the following:

$$l(\phi; y) = -\frac{1}{2} \left\{ N \log(2\pi) + \log |\mathbb{V}(\theta)| + (y - X\beta)^\top [\mathbb{V}(\theta)]^{-1} (y - X\beta) \right\}, \tag{2.5}$$

where $\mathbb{V}(\theta)$ is is the covariance matrix of $y$.

To reduce the bias in the process of estimating the components of the vector $\boldsymbol{\theta}$ by RMLE, [12] and [6] proposed to use a linear transformation of the type $\boldsymbol{y}^* = \boldsymbol{U}^\top \boldsymbol{y}$ with $\mathbb{E}\left(\boldsymbol{U}\boldsymbol{y}^*\right) = \boldsymbol{0}$. The considered the matrix $\boldsymbol{U}$, such that $\boldsymbol{U}^\top \boldsymbol{U} = \boldsymbol{I}_n$ and $\boldsymbol{U}\boldsymbol{U}^\top = \boldsymbol{I}_n - \boldsymbol{X}(\boldsymbol{X}^\top \boldsymbol{X})^{-1}\boldsymbol{X}^\top$, with

$$\boldsymbol{X} = \left(\boldsymbol{X}_i^T, ..., \boldsymbol{X}_n^T\right)^T,$$

so that,

$$\boldsymbol{y}^* \sim \mathcal{N}_{N-p}\left[\boldsymbol{0}, \boldsymbol{U}^\top \mathbb{V}(\boldsymbol{\theta})\boldsymbol{U}\right], \tag{2.6}$$

with

$$N = \sum_{i=1}^{n} m_i.$$

The logarithm of the restricted marginal likelihood function is as follows:

$$l_R(\boldsymbol{\theta}; \boldsymbol{y}) = -\frac{(N-p)}{2}\log(2\pi) - \frac{1}{2}\log|\mathbb{V}(\boldsymbol{\theta})| - \frac{1}{2}\log\left|\boldsymbol{X}^\top[\mathbb{V}(\boldsymbol{\theta})]^{-1}\boldsymbol{X}\right| - \frac{(N-p)}{2}\widehat{\boldsymbol{e}}^\top[\mathbb{V}(\boldsymbol{\theta})]^{-1}\widehat{\boldsymbol{e}}, \tag{2.7}$$

where $\widehat{\boldsymbol{e}} = \boldsymbol{y} - \boldsymbol{X}\widehat{\boldsymbol{\beta}}(\boldsymbol{\theta})$, and $\widehat{\boldsymbol{\beta}}(\boldsymbol{\theta}) = \left(\boldsymbol{X}^\top[\mathbb{V}(\boldsymbol{\theta})]^{-1}\boldsymbol{X}\right)^{-1}\boldsymbol{X}^\top[\mathbb{V}(\boldsymbol{\theta})]^{-1}\boldsymbol{y}$ is the MLE of $\boldsymbol{\beta}$ assuming that $\boldsymbol{\theta}$ is known. The maximization of (2.7) generates the estimators $\widehat{\boldsymbol{\theta}}_R$ of maximum plausibility of $\boldsymbol{\theta}$. Hence, $\widehat{\boldsymbol{\theta}}_R$ in conjunction with $\widehat{\boldsymbol{\beta}}_R = \widehat{\boldsymbol{\beta}}(\widehat{\boldsymbol{\theta}}_R)$ are the desired RMLEs [2].

The function (2.7) can also be written as follows:

$$l_R(\boldsymbol{\theta}; \boldsymbol{y}) = -\frac{1}{2}\left\{(N-p)\log(2\pi) + \log|\mathbb{V}(\boldsymbol{\theta})| + \log\left|\boldsymbol{X}^\top[\mathbb{V}(\boldsymbol{\theta})]^{-1}\boldsymbol{X}\right| + \boldsymbol{y}^\top \boldsymbol{P}\boldsymbol{y}\right\}, \tag{2.8}$$

with

$$\boldsymbol{P} = [\mathbb{V}(\boldsymbol{\theta})]^{-1} - [\mathbb{V}(\boldsymbol{\theta})]^{-1}\boldsymbol{X}\left(\boldsymbol{X}^\top[\mathbb{V}(\boldsymbol{\theta})]^{-1}\boldsymbol{X}\right)^{-1}\boldsymbol{X}^\top[\mathbb{V}(\boldsymbol{\theta})]^{-1}.$$

In general, the process that generates the vector $\boldsymbol{y}$ is not known (i.e., the true probability density $g(\boldsymbol{y})$ is not known). An LMEM is usually proposed assuming that the distribution of both, the random effects and the random error term are known. Here, we consider that $f(\boldsymbol{y}, \boldsymbol{\phi})$ is the density function of the random vector $\boldsymbol{y}$. If there exists a vector $\boldsymbol{\phi}_0 \in \boldsymbol{\Theta}$ such that $g(\boldsymbol{y}) = f(\boldsymbol{y}, \boldsymbol{\phi}_0)$, with $\boldsymbol{\Theta}$ being a compact subset of a $p$-dimensional Euclidean space, it can be concluded that the model would be correctly specified. Otherwise, the model would have a specification error. [20] illustrates that when the model is correctly specified, $\widehat{\boldsymbol{\phi}}_n$, which is obtained by either maximum likelihood or restricted maximum likelihood, is a consistent estimator for $\boldsymbol{\phi}_0$ [17, p.34], this is,

$$\widehat{\boldsymbol{\phi}}_n \xrightarrow{p} \boldsymbol{\phi}_0. \tag{2.9}$$

When the model is incorrectly specified, there exists a vector $\boldsymbol{\phi}^* \in \boldsymbol{\Theta}$, which minimizes the information criterion using the Kullback $-$ Leibler (KL) distance, that is,

$$KL(g : f, \boldsymbol{\phi}) = \mathbb{E}_g\left[\log\frac{g(\boldsymbol{y})}{f(\boldsymbol{y}, \boldsymbol{\phi})}\right] = \int_{\mathbb{R}^N} g(\boldsymbol{y})\log\frac{g(\boldsymbol{y})}{f(\boldsymbol{y}, \boldsymbol{\phi})}d\boldsymbol{y}. \tag{2.10}$$

However, when the model is properly specified, [19] proved that the only value that minimized the KL criterion was $\boldsymbol{\phi}^* = \boldsymbol{\phi}_0$.

Let us assume that the following matrices exist:

$$A(\boldsymbol{\phi}) = \mathbb{E}\left[\left(\frac{\partial^2 l(\boldsymbol{\phi}; y)}{\partial \phi_k \partial \phi_l}\right)\right], \qquad B(\boldsymbol{\phi}) = \mathbb{E}\left[\left(\frac{\partial l(\boldsymbol{\phi}; y)}{\partial \phi_k}\frac{\partial l(\boldsymbol{\phi}; y)}{\partial \phi_l}\right)\right],$$

$$A_n(\boldsymbol{\phi}) = \frac{1}{n}\sum_{i=1}^{n}\left(\frac{\partial^2 l_i(\boldsymbol{\phi}; y_i)}{\partial \phi_k \partial \phi_l}\right), \qquad B_n(\boldsymbol{\phi}) = \frac{1}{n}\sum_{i=1}^{n}\left(\frac{\partial l_i(\boldsymbol{\phi}; y_i)}{\partial \phi_k}\frac{\partial l_i(\boldsymbol{\phi}; y_i)}{\partial \phi_l}\right),$$

$$A(\boldsymbol{\phi}_0) = \mathbb{E}\left[\left(\frac{\partial^2 l(\boldsymbol{\phi}; y)}{\partial \phi_k \partial \phi_l}\right)\right]_{\boldsymbol{\phi}=\boldsymbol{\phi}_0}, \qquad B(\boldsymbol{\phi}_0) = \mathbb{E}\left[\left(\frac{\partial l(\boldsymbol{\phi}; y)}{\partial \phi_k}\frac{\partial l(\boldsymbol{\phi}; y)}{\partial \phi_l}\right)\right]_{\boldsymbol{\phi}=\boldsymbol{\phi}_0},$$

$$A_n(\widehat{\boldsymbol{\phi}}_n) = \frac{1}{n}\sum_{i=1}^{n}\left(\frac{\partial^2 l_i(\boldsymbol{\phi}; y_i)}{\partial \phi_k \partial \phi_l}\right)_{\boldsymbol{\phi}=\widehat{\boldsymbol{\phi}}_n}, \qquad B_n(\widehat{\boldsymbol{\phi}}_n) = \frac{1}{n}\sum_{i=1}^{n}\left(\frac{\partial l_i(\boldsymbol{\phi}; y_i)}{\partial \phi_k}\frac{\partial l_i(\boldsymbol{\phi}; y_i)}{\partial \phi_l}\right)_{\boldsymbol{\phi}=\widehat{\boldsymbol{\phi}}_n}.$$

According to [19],

$$A_n(\widehat{\boldsymbol{\phi}}_n) \xrightarrow{as} A(\boldsymbol{\phi}). \tag{2.11}$$

Now, let

$$-H = -\frac{\partial^2 l_i(\boldsymbol{\phi}; y_i)}{\partial \phi_k \partial \phi_l} \tag{2.12}$$

be the observed information matrix. If model (2.1) is correctly specified [19],

$$A(\boldsymbol{\phi}_0) + B(\boldsymbol{\phi}_0) = \mathbf{0}. \tag{2.13}$$

The expression (2.13) is called an equality of the information matrix. Under appropriate conditions [4], it is possible to demonstrate that

$$\sqrt{n}\widehat{\boldsymbol{\phi}}_n - \boldsymbol{\phi}_0 \xrightarrow{d} \mathcal{N}_s\left[\mathbf{0}, \mathbb{V}(\boldsymbol{\phi}_0)\right], \tag{2.14}$$

with

$$\mathbb{V}(\boldsymbol{\phi}_0) = [A(\boldsymbol{\phi}_0)]^{-1}B(\boldsymbol{\phi}_0)[A(\boldsymbol{\phi}_0)]^{-1}. \tag{2.15}$$

## 3. Proposed test

In this section, we describe in detail a new test designed to detect misspecification errors as and their source in Gaussian LMEMs. Our proposed test is based on the "Sandwich" estimator of the covariance matrix of $\widehat{\boldsymbol{\phi}}_n$ and in the equality of the information matrix given by (2.13), under the null hypothesis $H_0$ that the model (2.1) is correctly specified.

Considering the asymptotic distribution of the estimator $\widehat{\boldsymbol{\phi}}_n$ and under regularity conditions [19], it is known that

$$\sqrt{n}\left(\widehat{\boldsymbol{\phi}}_n - \boldsymbol{\phi}_0\right) \overset{A}{\sim} \mathcal{N}_s\left[\mathbf{0}, \mathbb{V}(\boldsymbol{\phi}_0)\right], \tag{3.1}$$

where $\mathbb{V}(\boldsymbol{\phi}_0)$ is as in (2.15).

Let us consider the vector

$$d(\boldsymbol{\phi}) = \operatorname{diag}\left\{[A(\boldsymbol{\phi})]^{-1}B(\boldsymbol{\phi})[A(\boldsymbol{\phi})]^{-1} + [A(\boldsymbol{\phi})]^{-1}\right\}, \tag{3.2}$$

and assume that it is differentiable in $\boldsymbol{\phi}_0$. Thus, we can build the following:

$$\boldsymbol{\nabla}d(\boldsymbol{\phi}_0) = \boldsymbol{\nabla}d(\boldsymbol{\phi})\big|_{\boldsymbol{\phi}=\boldsymbol{\phi}_0} \neq \mathbf{0}. \tag{3.3}$$

Equation (3.3) is a condition necessary for the existence of the gradient. It follows that the vector

$$\begin{aligned}
d_n^*(\widehat{\boldsymbol{\phi}}_n) &= \operatorname{diag}\left\{[A_n(\widehat{\boldsymbol{\phi}}_n)]^{-1}B_n(\widehat{\boldsymbol{\phi}}_n)[A_n(\widehat{\boldsymbol{\phi}}_n)]^{-1} + [A_n(\widehat{\boldsymbol{\phi}}_n)]^{-1}\right\} \\
&= \operatorname{diag}\left\{[A_n(\widehat{\boldsymbol{\phi}}_n)]^{-1}B_n(\widehat{\boldsymbol{\phi}}_n)[A_n(\widehat{\boldsymbol{\phi}}_n)]^{-1}\right\} + \operatorname{diag}\left\{[A_n(\widehat{\boldsymbol{\phi}}_n)]^{-1}\right\} \\
&= \frac{1}{n}\sum_{i=1}^{n} d_i^*(\boldsymbol{\phi})\big|_{\boldsymbol{\phi}=\hat{\boldsymbol{\phi}}_n}
\end{aligned} \tag{3.4}$$

is an estimator of (3.2), and a potential indicator for detecting specification errors in the model (2.1).

Let $d_n(\boldsymbol{\phi})$ be the following:

$$d_n(\boldsymbol{\phi}) = \operatorname{diag}\left\{[A_n(\boldsymbol{\phi})]^{-1}B_n(\boldsymbol{\phi})[A_n(\boldsymbol{\phi})]^{-1} + [A_n(\boldsymbol{\phi})]^{-1}\right\} = \begin{pmatrix} d_{11}(\boldsymbol{\phi}) \\ \vdots \\ d_{ss}(\boldsymbol{\phi}) \end{pmatrix}. \tag{3.5}$$

Then,

$$\boldsymbol{\nabla}d_n(\boldsymbol{\phi}) = \boldsymbol{\nabla}\left\{\operatorname{diag}\left[A_n(\boldsymbol{\phi})^{-1}B_n(\boldsymbol{\phi})[A_n(\boldsymbol{\phi})]^{-1} + [A_n(\boldsymbol{\phi})]^{-1}\right]\right\} \tag{3.6}$$

$$= \begin{pmatrix} \frac{\partial d_{11}(\boldsymbol{\phi})}{\partial \phi_1} & \cdots & \frac{\partial d_{11}(\boldsymbol{\phi})}{\partial \phi_s} \\ \vdots & \ddots & \vdots \\ \frac{\partial d_{ss}(\boldsymbol{\phi})}{\partial \phi_1} & \cdots & \frac{\partial d_{ss}(\boldsymbol{\phi})}{\partial \phi_s} \end{pmatrix}_{s\times s}. \tag{3.7}$$

Using (3.1) under $H_0$ and the Delta method [17, p.136], we obtain the following:

$$\sqrt{n}\left(d_n^*(\widehat{\boldsymbol{\phi}}_n) - d(\boldsymbol{\phi}_0)\right) \xrightarrow{D} \mathcal{N}_s(\mathbf{0}, \boldsymbol{\nabla}d(\boldsymbol{\phi}_0)\mathbb{V}(\boldsymbol{\phi}_0)\boldsymbol{\nabla}d(\boldsymbol{\phi}_0)^\top), \qquad n \to \infty. \tag{3.8}$$

Now, if we consider (3.8) under $H_0$ and use Cochran's theorem [17, p.137], then the test statistic of the alternative "Sandwich" estimator (ASEST) takes the following form

$$\text{ASEST} = n d_n^*(\widehat{\boldsymbol{\phi}}_n)^\top \left[\boldsymbol{\nabla}d_n(\widehat{\boldsymbol{\phi}}_n)\widehat{\mathbb{V}}(\widehat{\boldsymbol{\phi}}_n)\boldsymbol{\nabla}d_n(\widehat{\boldsymbol{\phi}}_n)^\top\right]^{-1} d_n^*(\widehat{\boldsymbol{\phi}}_n) \tag{3.9}$$

where

$$\widehat{\mathbb{V}}(\widehat{\boldsymbol{\phi}}_n) = [A_n(\widehat{\boldsymbol{\phi}}_n)]^{-1}B_n(\widehat{\boldsymbol{\phi}}_n)[A_n(\widehat{\boldsymbol{\phi}}_n)]^{-1} \tag{3.10}$$

is an unbiased and consistent estimator of the covariance matrix of the "Sandwich" estimator for $\widehat{\boldsymbol{\phi}}$ [19]. It is straightforward to show that, under $H_0$, $\text{ASEST} \sim \chi_s^2$ as $n \to \infty$. Thus, the Type I error of the test can be calculated as follows:

$$\text{Type I error} = P(\text{ASEST} > \chi^2_{\alpha,s} \,|\, H_0). \tag{3.11}$$

Similarly, under the alternative hypothesis $H_1$, the power of the test can be calculated as

$$\text{Power} = 1 - P(\text{ASEST} < \chi^2_{\alpha,s} \,|\, H_1). \tag{3.12}$$

## 4. Simulation study

We conducted numerical experiments to study the behavior of our proposed test in terms of the Type I error rate and the Power to identify potential misspecifications in LMEMs. The structure in all models is the same as in (2.1).

In order to reflect particular situations we commonly encounter in practical contexts, five different cases were considered:

- *Case I:* Misspecification of the systematic component;
- *Case II:* Misspecification of the number of random effects; and
- *Case III:* Random effects are considered independent;
- *Case IVa:* Misspecification of the within-subject covariance structure;
- *Case IVb:* Misspecification of the within-subject covariance structure.

We simulated data following a previously published structure of a Gaussian LMEM [1], and considered [15] to identify a correctly specified model. In particular, the individuals' profiles were used to identify both the structure of the mean response and of the random effects of the LMEM. The identified model, shown as *Case I* in Table 1, is a second-degree polynomial in the fixed part, and a first-degree polynomial in the random part. This model can be written as follows:

$$y_{ij} = \beta_0 + \beta_1 x_{ij} + \beta_2 x_{ij}^2 + b_{0i} + b_{1i} x_{ij} + e_{ij}, \tag{4.1}$$

where $y_{ij}$ represents a $j$-th observation of the $i$-th the individual, $x_{ij}$ is the $j$-th time registered for the $i$-th individual, $\boldsymbol{\beta} = (\beta_0, \beta_1, \beta_2)$ are the location parameters, $\boldsymbol{\theta} = (b_{0i}, b_{1i})$ represents the intercept, and random slope, and $e_{ij}$ is the random error term.

Model (4.1) can be written in a compact form where $\boldsymbol{X}_i = \begin{pmatrix} 1 & x_{i_1} & x_{i_1}^2 \\ \vdots & \vdots & \vdots \\ 1 & x_{i_{m_i}} & x_{i_{m_i}}^2 \end{pmatrix}$ and $\boldsymbol{Z}_i = \begin{pmatrix} 1 & x_{i_1} \\ \vdots & \vdots \\ 1 & x_{i_{m_i}} \end{pmatrix}$.

Assuming that $\boldsymbol{G} = \begin{pmatrix} \sigma_0^2 & \sigma_{01} \\ \sigma_{01} & \sigma_1^2 \end{pmatrix}$, $\boldsymbol{R}_i = \sigma^2 \boldsymbol{I}_n$, $\boldsymbol{b}_i \sim \mathcal{N}_2[\boldsymbol{0}, \boldsymbol{G}]$, $\boldsymbol{e}_i \sim \mathcal{N}_{m_i}[\boldsymbol{0}, \boldsymbol{R}_i]$, and that $\boldsymbol{b}_i$ and $\boldsymbol{e}_i$ are independent $(i = 1, \ldots n)$, it follows that

$$\boldsymbol{y}_i \sim \mathcal{N}_{m_i}(\boldsymbol{X}_i \boldsymbol{\beta}, \mathbb{V}_i), \qquad \mathbb{V}_i = \boldsymbol{Z}_i \boldsymbol{G} \boldsymbol{Z}_i^\top + \boldsymbol{R}_i. \tag{4.2}$$

Let $\boldsymbol{\phi} = (\boldsymbol{\beta}^\top, \boldsymbol{\theta}^\top)^\top$ be the vector of all model parameters for model (4.1). Using LME, the parameter estimates $\widehat{\boldsymbol{\phi}}_n$ of model (4.1) are as follows:

$$\widehat{\boldsymbol{\beta}} = \begin{pmatrix} 2.451 & 1.390 & 0.548 \end{pmatrix}^\top, \ \widehat{\boldsymbol{G}} = \begin{pmatrix} 0.536 & 0.569 \\ 0.569 & 0.888 \end{pmatrix}, \text{ and } \widehat{\sigma}^2 = 0.404, \tag{4.3}$$

where $\widehat{\boldsymbol{G}}$ is the estimated covariance matrix of the vector of random effects, $\boldsymbol{G}$. These values are considered to be the true values of the parameters in the simulation process. Hence, we subsequently simulate $B$ data sets of size $n$ and fit a model according to Table 1. The next step is to calculate the power and the Type I error rate of the proposed test. For the former, we simulated data from the identified model (that is, *Case I*) and fit a model in accordance with Table 1. Similarly, to calculate the Type I error rate, we simulate data of the identified model and fit a different model.

**Table 1.** Structure, fixed effects, random effects, and error term considered in the numerical experiments.

| Case | Model | $G$ | $R_i$ |
|---|---|---|---|
| I | $y_{ij} = \underbrace{\beta_0 + \beta_1 x_{ij} + \beta_2 x_{ij}^2}_{\text{Fixed effects}} + \underbrace{b_{0i} + b_{1i} x_{ij}}_{\text{Random effects}} + \underbrace{e_{ij}}_{\text{Error term}}$ | $\begin{pmatrix} \sigma_0^2 & \sigma_{01} \\ \sigma_{01} & \sigma_1^2 \end{pmatrix}$ | $\sigma^2 I_n$ |
| II | $y_{ij} = \underbrace{\beta_0 + \beta_1 x_{ij}}_{\text{Fixed effects}} + \underbrace{b_{0i} + b_{1i} x_{ij}}_{\text{Random effects}} + \underbrace{e_{ij}}_{\text{Error term}}$ | $\begin{pmatrix} \sigma_0^2 & \sigma_{01} \\ \sigma_{01} & \sigma_1^2 \end{pmatrix}$ | $\sigma^2 I_n$ |
| III | $y_{ij} = \underbrace{\beta_0 + \beta_1 x_{ij}}_{\text{Fixed effects}} + \underbrace{b_{0i}}_{\text{Random effect}} + \underbrace{e_{ij}}_{\text{Error term}}$ | $\sigma_0^2$ | $\sigma^2 I_n$ |
| IVa | As in Case I | $\begin{pmatrix} \sigma_0^2 & \sigma_{01} \\ \sigma_{01} & \sigma_1^2 \end{pmatrix}$ | $\sigma^2 \begin{pmatrix} 1 & \rho & \rho^2 & \cdots & \rho^{m_i-1} \\ & 1 & \rho & \cdots & \rho^{m_i-2} \\ & & 1 & \rho & \vdots \\ & & & \ddots & \rho \\ & & & & 1 \end{pmatrix}, |\rho| < 1$ |
| IVb | As in Case I | $\begin{pmatrix} \sigma_0^2 & 0 \\ 0 & \sigma_1^2 \end{pmatrix}$ | $\sigma^2 I_n$ |
| VI | $y_{ij} = \underbrace{\beta_0 + \beta_1 x_{ij} + \beta_2 x_{ij}^2}_{\text{Fixed effects}} + \underbrace{b_{0i}}_{\text{Random effect}} + \underbrace{e_{ij}}_{\text{Error term}}$ | $\sigma_0^2$ | $\sigma^2 \begin{pmatrix} 1 & \rho & \rho^2 & \cdots & \rho^{m_i-1} \\ & 1 & \rho & \cdots & \rho^{m_i-2} \\ & & 1 & \rho & \vdots \\ & & & \ddots & \rho \\ & & & & 1 \end{pmatrix}, |\rho| < 1$ |

## 5. Results

### 5.1. Case I: Misspecification of the systematic component

Here, we generated $B = 10{,}000$ data sets from Model *I* in Table 1 using the estimates $\widehat{\boldsymbol{\beta}}, \widehat{\boldsymbol{G}}$, and $\widehat{\sigma}^2$ in (4.3). Thus, the response vector is such that $\boldsymbol{y}_i \sim \mathcal{N}_{m_i}(\boldsymbol{X}_i\boldsymbol{\beta}, \mathbb{V}_i)$, where $\mathbb{V}_i = \boldsymbol{Z}_i\boldsymbol{G}\boldsymbol{Z}_i^\top + \boldsymbol{R}_i$.

In order to assess the misspecification of the systematic component, we fitted Model *II* in Table 1 to the simulated data. Note that, in this model, the quadratic term is omitted; hence, the model is misspecified. The Type I error rate and Power, as a function of $n$, are reported in Figure 1. Overall, our results for Case *I* indicate that as $n \to \infty$, the Type I error of the test gets closer to the nominal Type I error rate of 5% and the Power increases, suggesting that our proposed test is capable of correctly detecting a misspecificacion in the systematic component of the LMEM model as $n$ increases.

## 5.2. Case II: Misspecification of the number of random effects

To assess the Type I error and Power of the proposed test in this specific instance, we created $B = 10,000$ datasets based on Model *I* of Table 1, utilizing the estimates provided in (4.3). Subsequently, we fitted Model *III* from the same table. By design, here we omit the random slope that induces the specification error in the number of random effects; therefore, the covariance matrix generated by this model induces a uniform structure, in which the variances are constant over time. Our results, presented in Figure 1, show that as $n \to \infty$, the Type I error of the test gets closer to the nominal level and the Power increases, suggesting that the proposed test is capable of correctly detecting the misspecification in the number of the random effects as the sample size increases.

## 5.3. Case III: Random effects are considered independent

Using a simulation strategy similar to that previously discussed in Case *I* and Case *II*, we simulated $B = 10,000$ data sets from Model *I* and fitted Model *III* (see Table 1 for more details) to reflect a misspecification error in the covariance matrix of the random effects. In particular, we wrongly considered that the random effects were independent and assessed the performance of the proposed test to detect such misspecifications. Figure 1 shows the results. Overall, our findings suggest that the Type I error of the proposed test gets closer to the nominal level of 5% as $n$ increases, suggesting that our test controls the probability of wrongly identifying a misspecification of the covariance structure of the random effects. In addition, the Power of the proposed test in Case *III* is $> 0.95$ regardless of $n$, and increases to 1 for $n > 300$. This result indicates that the proposed test is highly likely to detect a misspecification error in the covariance structure of the random effects when it actually exists.

## 5.4. Case IVa: Misspecification of the within-subject covariance structure

Here, we assess the performance of the proposed test when the covariance matrix of the response vector is misspecified. In particular, we consider that the random error term follows an autoregressive (AR) process of order 1 with $\rho = 0.9$. Following (3.9), the test statistic is $\chi^2_{\text{calculated}} = 1233.4$ and the associated $p$-value is $< 1 \times 10^{-16}$, indicating that the LMEM is misspecified.

Following a similar simulation strategy to that previously described, we assessed the Type I error rate and the Power of the proposed test when a misspecification of the within-subject covariance structure exists. The results are presented in Figure 1. As seen for the other cases, for Case *IVa*, the Type I error of the proposed test gets closer to the nominal level as $n$ increases, and the Power of the proposed test is $> 0.9$ regardless of $n$ and increases to 1 for $n > 500$. Overall, these results indicate that the proposed test performs reasonably well for detecting a misspecification error in the within-subject covariance structure when it actually exists, and controls the Type I error when it does not.
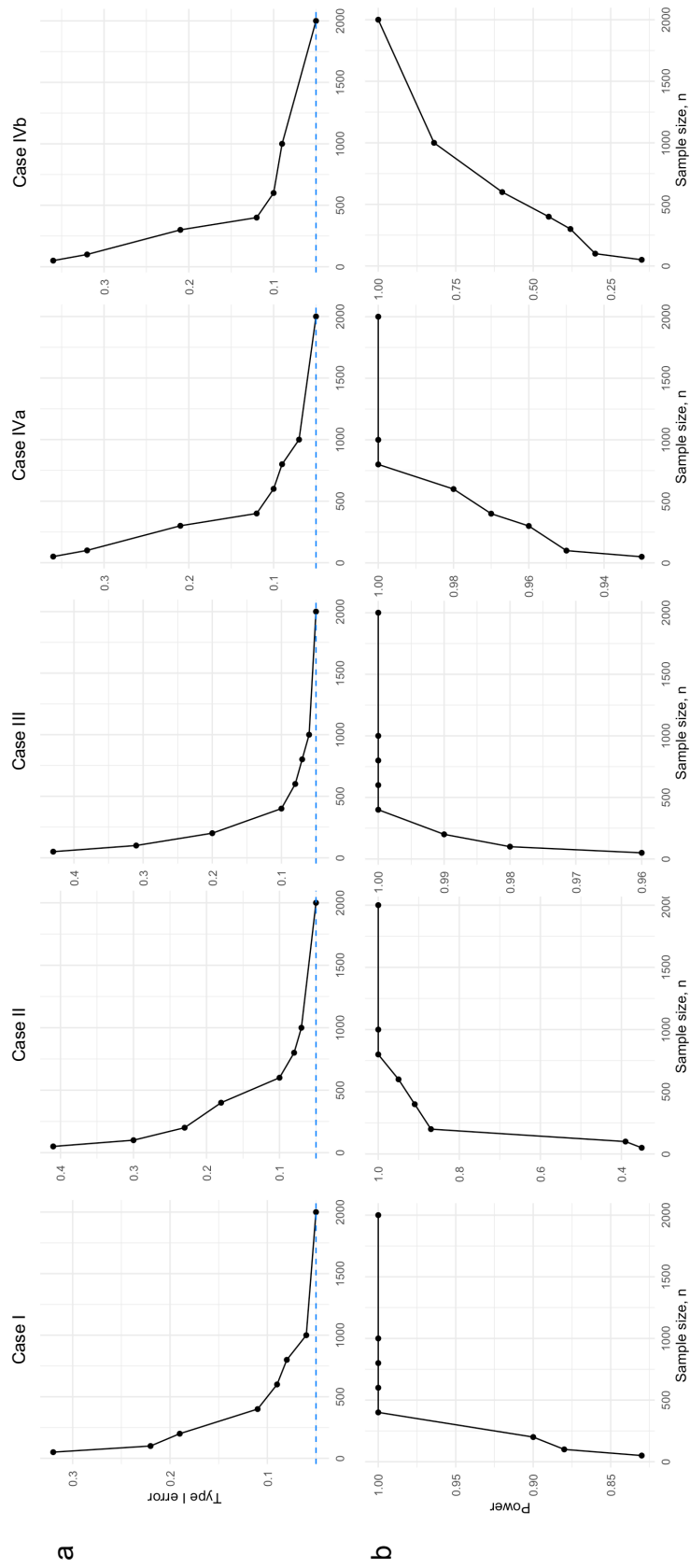
**Figure 1. (a)** Type I error and **(b)** Power of the proposed test for each Case in Table 1 as a function of the sample size, *n*. The blue line represents a Type I error of 5%.

*5.5. Case IVb: Misspecification in the within-subject covariance structure*

This case is similar to Case *IVa*. However, there is a specification error in the number of random effects and in the within-subject covariance matrix. In particular, we simulated $B = 10,000$ data sets from a model that has an intercept, slope, and quadratic effect in the fixed part, and only an intercept in the random part. In addition, the random error term follows an AR(1) process. With the simulated data, we fitted a Model *I* of 1, which includes an intercept, slope, and square effect in the fixed part, and an intercept and slope in the random part. We subsequently fitted Model *IVb*, which includes an intercept and a slope in the fixed part, and considers that the random effects are homoscedastic with a conditional independence. Figure 1 displays the results.

## 6. An illustration with real data

Nagle (2018) [11] presented a dataset that analyzed the temporal shift in the production of stop consonants by a group of 24 English learners (Figure 2). The phonetic context is controlled with 4 dummy characters: 'Pafo', 'Bafo', 'Pamuso', and 'Bamuso'. With the first two characters, the occlusion occurs on a stressed syllable, while with the other two, the occlusion occurs on an unstressed syllable. The outcome variable was voice onset time (VOT), which is an acoustic measure representing the time elapsed between the onset of vocalisation and the release of an occlusion closure. Five sessions were conducted for each student. Two stress categories were analysed, as well as each participant's age and stress, among other variables. In this study, it was of interest in the study to differentiate the VOT variation from the individuals' variations.
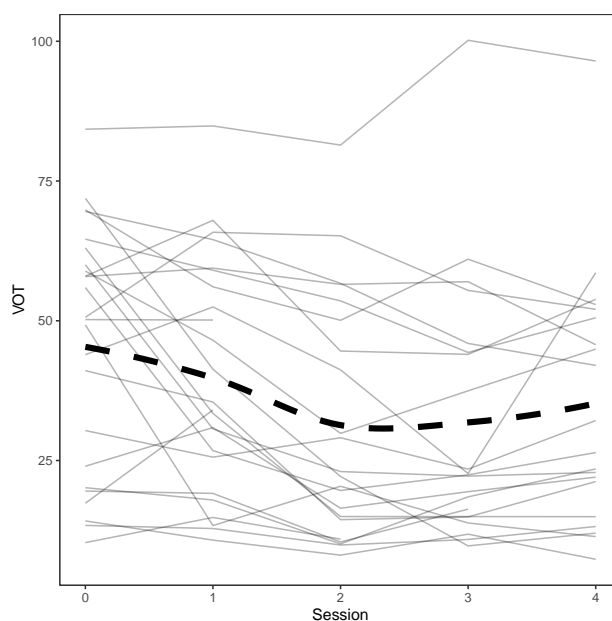


**Figure 2.** Individuals' profiles (in gray) around the loess profile (dashed line).

Figure 2 suggests that introducing sessions as a random effect may be appropriate. Thus, an LMEM will allow us to analyze the individuals' and the global behavior, even with incomplete data for some individuals, as shown in Figure 2. In addition, we illustrate the proposed test to identify the source of error, if there is a misspecification in the model.

## 6.1. *Model specification, estimation and proposed test*

Given the behavior of the profiles in Figure 2, we can suggest an LMEMs with intercept, slope, and quadratic effect in the fixed part, and an intercept and slope in the random part. To model the observed heteroscedasticity, we propose a homoscedastic LMEM with a conditional independence for simplicity. The general form of this model was previously described and is given in (4.1). Table 2 shows the MLEs of $\phi = (\beta^{\top}, \theta^{\top})^{\top}$, and the vector of all model parameters in (4.1).

**Table 2.** MLEs for the (a) Fixed and (b) Random part of model (4.1).

(a) Fixed part

| Term | Estimate | SE | *df* | *t* | *P*−value |
|------|----------|------|---------|-------|-----------|
| $\beta_0$ | 46.24 | 4.16 | 24.22 | 11.12 | 0.00 |
| $\beta_1$ | -10.22 | 1.33 | 28.69 | -7.66 | 0.00 |
| $\beta_2$ | 1.98 | 0.16 | 4356.48 | 12.34 | 0.00 |

(b) Random part

| Term | Variance Estimate |
|------|-------------------|
| $b_0$ | 405.96 |
| $b_1$ | 33.55 |
| $e$ | 15.984 |

Note: SE: Standard Error; *df*: degrees of freedom; *t*: test statistic.

Under the null hypothesis that the model (4.1) is correctly specified, the ASEST test statistic is $\chi^2_{\text{calculated}} = 1.279 \times 10^{-9}$ and the associated *p*-value is $> 0.05$. Thus, no misspecification is detected.
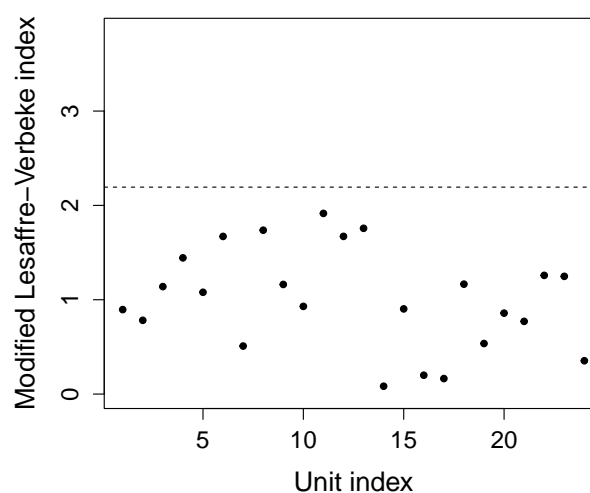


**Figure 3.** Modified Lesaffre–Verbeke unit index plot for model (4.1). Dashed line represents the $Q_3 + 1.5$ *IQR*. $Q_3$: 3rd quatile; *IQR*: interquartile range.

## 6.2. *Graphical diagnostic tests*

As a complement to the ASEST proposed test for identifying misspecification errors in model (4.1) for the VOT data set [11], here we employed several graphical diagnostic tests on the fitted model. Our

results are presented in Figures 3–8.

Figure 3 presents the residual diagnostic plot for the modified Lesaffre–Verbeke index. The plot suggests that the proposed covariance structure may be suitable for the 24 units, as our results do not indicate any significant deviations from the expected behavior if the model were incorrectly specified.

Our interest in this section is to illustrate the graphical diagnostic methodology that allows us to identify and locate the source of the specification error, if it exists. In this case, the model is incorrectly specified in the sources identified.

On the other hand, Figure 4 shows the standard raw residual Q-Q plot for the fitted model. This result, which is in line with the findings of the proposed test, suggests that the error distribution for the adopted LMEMs does not seem to have heavy tails.
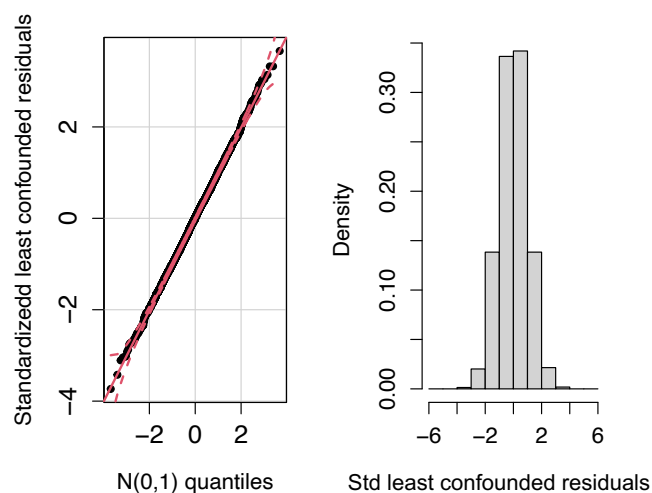


**Figure 4.** Q-Q plot (*left*) and histogram (*right*) for the standardized conditional residuals.

Additionally, we analyzed the potential influential observations based on the Mahalanobis distance (Figure 5). Notably, it is critical to identify observations may be crucial for a subsequent analysis in the context LMEMs; however, this cannot be achieved using our ASEST test. Therefore, the relevance of this plot lies in its ability to visualize and potentially highlight influential observations. Interestingly, we identified that two possible observations (i.e., observations #5 and #6) may be influential. Although this result warrants a further investigation to detect potential outliers or inconsistent values based on the findings, we can statistically treat them as possible influential values.

Another important aspect that cannot be addressed using the analytical version of the ASEST test, is the identification of patterns in the fixed effects of the LMEM. Figure 6 shows the diagnostic plot for such a suggestion when model (4.1) is fitted to the VOT data (Figure 2), which suggests that there is no evidence to illustrate the omission of any fixed effect.
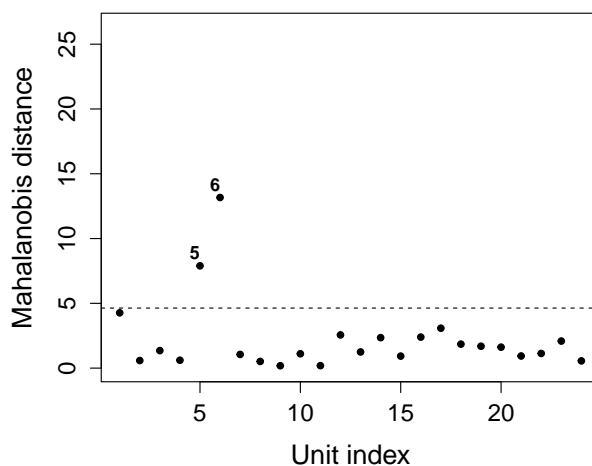
**Figure 5.** Mahalanobis distance as a function of the unit index for identifying influential observations.
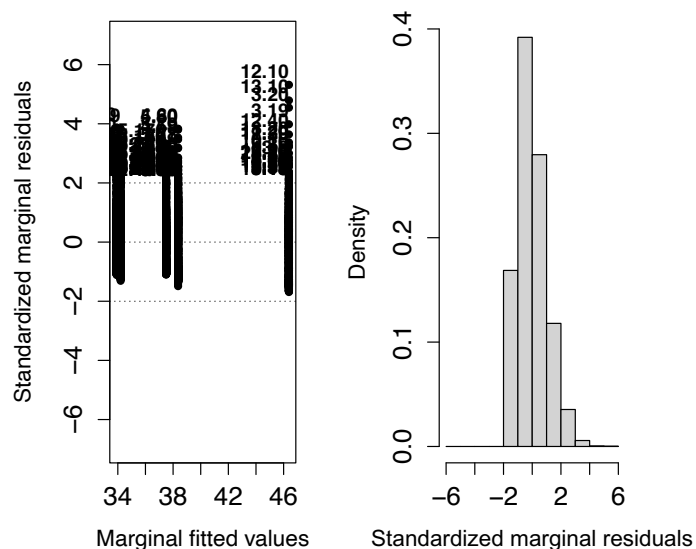


**Figure 6.** Standardized vs. marginal fitted residuals (*left*) and histogram of the standardized residuals (*right*) for the VOT data set based on model (4.1).

Consequently, we should check whether the adopted Gaussian assumption for the random effect in model (4.1) is not suitable through the ASEST test. Thus, we propose the use of a Q-Q plot for the Mahalanobis distance, as shown in Figure 7. Overall, these results suggest no evidence against the adopted Gaussian assumption for the random effects.

**Figure 7.** Q-Q plot for the Mahalanobis distance.

Finally, the vector of conditional errors was assumed to have a homoscedastic structure in model (4.1). However, the ASEST test cannot be directly used to validate such an assumption. Thus, we encourage the use of the standardized minimally confounded residuals, shown in Figure 8, as a suitable alternative for such a purpose. When applied to the VOT data set, the plot indicates that the assumption holds, which is in line with the assumption made while proposing the LMEM model given in (4.1).



**Figure 8.** Standardized residuals vs. predicted values (*left*) and histogram of the standardized residuals (*right*) for the VOT data set based on model (4.1).

## 7. Discussion and conclusions

LMEMs are widely used to analyze complex data structures; however, they are susceptible to various types of misspecification errors that can lead to biased or inconsistent estimates. Hence, detecting misspecification errors in LMEMs is crucial to ensure the accuracy and reliability of statistical inferences. These errors can arise from incorrect assumptions about the systematic component, random effects, within-subject covariance structure, or error term covariance structure.

In this study, we developed ASEST, a novel test specifically designed to identify misspecifications within LMEMs and accurately identifying their sources. Our approach applied the Delta method to leverage the asymptotic behavior of the test. The construction methodology was outlined in Sections 2 and 3, followed by a comprehensive analysis of diverse scenarios explained in detail in Section 4, and the associated results in Section 5.

One of the main advantages of our proposed test is its computational efficiency. We developed and implemented functions that exhibited rapid responsiveness, which ensures their applicability for real-world use. In the future, we plan to contribute the implementation of the ASEST test to the Comprehensive R Archive Network (CRAN) [13] repository to facilitate its accessibility and wider use.

An intriguing aspect of our methodology lies in its compatibility with models that feature no missing data, as displayed in Figure 2. Furthermore, we complement our analytical framework with a graphical methodology. This combined approach not only discerns erroneous model specifications, but also visually elucidates the origins of misspecification within the model structure. Notably, few methodologies are specifically designed to achieve this goal, particularly when it is of interest to identify the source of the misspecification. However, the integration of a robust statistical test, methodological clarity, computational efficiency, and the potential for future repository integration underscore the substantive contributions of our study. Our approach holds promise for a widespread adoption, offering a comprehensive toolkit to accurately assess and rectifying misspecification errors within LMEMs. Hence, the integration of statistical testing and graphical diagnostics in our methodology will significantly expand the diagnostic capabilities of LMEMs.

Future studies should be aimed at investigating additional scenarios to first assess the impact of simultaneous changes in the distribution of random effects and random errors. In this complex setting, it is crucial to analyse its feasibility to simultaneously detect the origin of both errors. This is particularly important as it is challenging to distinguish between the two. Identifying the sources of these errors in this scenario poses significant computational challenges, which require innovative methods and advanced statistical techniques to effectively address this issue.

Second, the performance of our test must be assessed in cases where the vector of random effects does not follow a Normal distribution. Although preliminary numerical experiments showed that the results presented in this study seemed to hold when this was the case, a more comprehensive evaluation has yet to be explored and completed. Third, here we considered a homoscedastic conditional independence structure for the covariance of the random error term. Further lines of research could explore $AR(p)$, non-constant, and not homogeneous covariance structures for the error term.

Finally, in our numerical experiments (see Section 4 for further details), the number of predictors or features $p$ were small compared to the sample size $n$. However, in many fields, the number of features

in the data are typically larger than the sample size (i.e., $p >>> n$). Therefore, the performance of our proposed test must be evaluated in these settings to ensure its effectiveness in real-world applications.

## Code availability

The R code for generating the plots and results in this paper is available from the first author upon request.

## Author contributions

Jairo A. Angel: Conceptualization, Methodology, Software, Validation, Formal analysis, Investigation, Resources, Data curation, Writing—original draft, Writing—review and editing, Visualization, Funding acquisition; Francisco M. M. Rocha: Conceptualization, Methodology, Software, Validation, Writing—original draft, Supervision, Project administration, Funding acquisition; Jorge I. Vélez: Formal analysis, Resources, Writing—review and editing, Visualization, Funding acquisition; Julio M. Singer: Conceptualization, Methodology, Software, Validation, Resources, Writing—original draft, Supervision, Project administration, Funding acquisition. All authors have read and agreed to the published version of the manuscript.

## Acknowledgements

## Conflict of interest

The authors declare that they have no conflict of interest in this work.

## References

1. J. Afiune, Avaliação Ecocardiografica Evolutiva de Recém−nascidos Pre-termo, do Nascimento Até o Termo, Phd thesis, Universidade de São Paulo, 2000.

2. E. Demidenko, *Mixed Models: Theory and Applications with R*, 2 Eds., Hoboken: John Wiley & Sons, Inc., 2013.

3. R. Drikvandi, G. Verbeke, G. Molenberghs, Diagnosing Misspecification of the Random-Effects Distribution in Mixed Models, *Biometrics*, **73** (2017), 63–71.

4. D. A. Freedman, On the so-called 'Huber Sandwich Estimator' and 'Robust Standard Errors', *Amer. Stat.*, **60** (2006), 299–302. http://doi.org/10.1198/000313006X152207

5. S. K. Hanneman, Design, Analysis, and Interpretation of Method-Comparison Studies, *AACN Adv. Crit. Care*, **19** (2008), 223–234. http://doi.org/10.4037/15597768-2008-2015

6. D. A. Harville, Maximum Likelihood Approaches to Variance Component Estimation and to Related Problems, *J. Amer. Stat. Assoc.*, **72** (1977), 320–338.

7. X. Huang, Detecting Random-Effects Model Misspecification via Coarsened Data, *Comput. Stat. Data Anal.*, **55** (2011), 703–714. https://doi.org/10.1016/j.csda.2010.06.012

8. J. Jiang, Goodness-of-Fit Tests for Mixed Model Diagnostics, *Ann. Stat.*, **29** (2001), 1137–1164.

9. N. Lange, L. Ryan, Assessing Normality in Random Effects Models, *Ann. Stat.*, **17** (1989), 624–642.

10. S. Litière, A. Alonso, G. Molenberghs, Type I and Type II Error under Random-Effects Misspecification in Generalized Linear Mixed Models, *Biometrics*, **63** (2007), 1038–1044.

11. C. Nagle, An Introduction to Fitting and Evaluating Mixed-effects Models in R, *Proceedings of the 10th Pronunciation in Second Language Learning and Teaching Conference*, 2018, 82–105.

12. H. D. Patterson, R. Thompson, Recovery of Inter-block Information when Block Sizes are Unequal, *Biometrika*, **58** (1971), 545–554,

13. R Core Team, R: A Language and Environment for Statistical Computing, Available from: `https://www.r-project.org/`.

14. G. K. Robinson, That BLUP is a Good Thing: The Estimation of Random Effects, *Stat. Sci.*, **6** (1991), 15–32,

15. F. M. Rocha, J. M. Singer, Selection of Terms in Random Coefficient Regression Models, *J. Appl. Stat.*, **45** (2018), 225–242.

16. H. Schielzeth, N. J. Dingemanse, S. Nakagawa, D. F. Westneat, H. Allegue, C. Teplitsky, et al., Robustness of Linear Mixed-Effects Models to Violations of Distributional Assumptions, *Methods Ecol. Evol.*, **11** (2020), 1141–1152. https://doi.org/10.1111/2041-210X.13434

17. P. Sen, J. Singer, *Large Sample Methods in Statistics: An Introduction with Applications*, Boca Raton: CRC Press, 1993.

18. J. M. Singer, F. M. Rocha, J. S. Nobre, Graphical Tools for Detecting Departures from Linear Mixed Model Assumptions and Some Remedial Measures, *Int. Stat. Rev.*, **85** (2017), 290–324. https://doi.org/10.1111/insr.12178

19. H. White, Maximum Likelihood Estimation of Misspecified Models, *Econometrica*, **50** (1982), 1–25. https://doi.org/10.2307/1912526

20. D. Yu, X. Zhang, K. K. Yau, Asymptotic Properties and Information Criteria for Misspecified Generalized Linear Mixed Models, *J. R. Stat. Soc. Ser. B (Stat. Methodol.)*, **80** (2018), 817–836.

AIMS Press