



Robust estimation of fixed effect parameters and variances of linear mixed models: the minimum density power divergence approach

Giovanni Saraceno^{1,3} · Abhik Ghosh² · Ayanendranath Basu² · Claudio Agostinelli¹

Received: 17 May 2022 / Accepted: 2 March 2023 / Published online: 29 March 2023
© The Author(s) 2023

Abstract

Many real-life data sets can be analyzed using linear mixed models (LMMs). Since these are ordinarily based on normality assumptions, under small deviations from the model the inference can be highly unstable when the associated parameters are estimated by classical methods. On the other hand, the density power divergence (DPD) family, which measures the discrepancy between two probability density functions, has been successfully used to build robust estimators with high stability associated with minimal loss in efficiency. Here, we develop the minimum DPD estimator (MDPDE) for independent but non-identically distributed observations for LMMs according to the variance components model. We prove that the theoretical properties hold, including consistency and asymptotic normality of the estimators. The influence function and sensitivity measures are computed to explore the robustness properties. As a data-based choice of the MDPDE tuning parameter α is very important, we propose two candidates as “optimal” choices, where optimality is in the sense of choosing the strongest downweighting that is necessary for the particular data set. We conduct a simulation study comparing the proposed MDPDE, for different values of α , with S-estimators, M-estimators and the classical maximum likelihood estimator, considering different levels of contamination. Finally, we illustrate the performance of our proposal on a real-data example.

Keywords Linear mixed models · Minimum density power divergence estimator · Robustness

Mathematics Subject Classification 62G35 · 62G05

✉ Giovanni Saraceno
gsaracen@buffalo.edu

Extended author information available on the last page of the article

1 Introduction

A major interest in statistics concerns the estimation of averages and their variation. The most commonly used method for this purpose is, probably, the *Linear Model* (LM). In this model, to give an example from a two-way layout, the expected value (mean) μ_{ij} of an observation y_{ij} may be expressed as a linear combination of unknown parameters such as $\mu_{ij} = \mu + \alpha_i + \beta_j$, where μ , α_i and β_j are the constants which we are interested in estimating. The linearity in the parameters means that we can write a linear model in the form $y_i = X_i\beta + \epsilon_i$, where β is the vector of unknown parameters, and the X_i s are known matrices. This formulation is the same as that used in case of linear regression model. In the present work, we consider the linear mixed models (LMMs), in which some (unknown) parameters are not treated as constants but as random variables. Random terms come into play when some items cannot be considered as fixed quantities, although their distributions are of interest. Hence, they are the tools to generalize the results to the entire population under study. The types of data that may be appropriately analyzed by LMMs include (i) *Clustered data*, where the dependent variable is measured once for each subject (the unit of analysis) and the units of analysis are grouped into, or nested within, clusters; (ii) *Repeated-measures data*, where the dependent variable is measured more than once on the same unit of analysis across levels of a factor, which may be time or other experimental conditions; (iii) *Longitudinal data*, where the dependent variable is measured at several points in time for each unit of analysis. For a general review of LMs and LMMs, see McCulloch and Searle (2001).

The standard methods used to estimate the parameters in LMMs are methods of maximum likelihood and restricted maximum likelihood. Generally, LMMs are based on normality assumptions, and it is well-known that these classical methods are not robust and can be greatly affected by the presence of small deviations from the assumptions. Furthermore, outlier detection for modern large data sets can be very challenging and, in any case, robust techniques cannot be replaced by the application of classical methods on outlier deleted data.

To answer the need for robust estimation in linear mixed models, a few methods have been proposed. The initial attempts were based on weighted versions of the log-likelihood function (see Huggins 1993a, b; Huggins and Staudte 1994; Stahel and Welsh 1994; Richardson and Welsh 1995; Richardson 1997; Welsh and Richardson 1997). Another attempt, discussed in Welsh and Richardson (1997), of robustifying linear mixed models consists of replacing the Gaussian distribution by the Student's t distribution (see also Lange et al. 1989; Pinheiro et al. 2001). However, this modification of the error distribution is intractable and complicated to implement. In Copt and Victoria-Feser (2006), a multivariate high breakdown point S-estimator, namely the CVFS-estimator, has been adapted to the linear mixed models setup, while the estimator given by Koller (2013), namely the SMDM-estimator, attempts to achieve robustness by a robustification of the score equations. Robust estimators have been proposed, more generally, for generalized linear mixed models by Yau and Kuk (2002) and Sinha (2004).

The density power divergence (DPD) (Basu et al. (1998)), which measures the discrepancy between two probability density functions, has been successfully used to build robust estimators for independent and identically distributed observations. In Ghosh and Basu (2013), the construction of the DPD and the corresponding minimum DPD estimator (MDPDE) has been extended to the case of independent but non-identically distributed data. This approach and theory covers the linear regression model and has later been extended to more general parametric regression models (Ghosh and Basu 2016, 2019; Castilla et al. 2018, 2021; Ghosh 2019, etc.). This MDPDE has become widely popular in recent times due to its good (asymptotic) efficiency along with high robustness, easy computability and direct interpretation as an intuitive generalization of the maximum likelihood estimator (MLE).

In the present work, our aim is to propose a robust estimator of the fixed effect parameters and variances of random effects under the linear mixed model set up. In particular, we are going to consider independent random effects according to the standard variance components models. This is done by an application of the definition and properties of the MDPDE, as formulated by Ghosh and Basu (2013), to the linear mixed models scenario. However, in this case, the observations have a complicated, non-identically distributed structure, which makes this adaptation non-trivial, in the computation of the estimator as well as in the theoretical derivations. We show that under appropriate conditions on the model matrices, the asymptotic and robustness properties hold, and the resulting estimator outperforms the competing robust estimators both in the presence and in the absence of contamination in sample data. Furthermore, the calculation of the sensitivity measures allows us to find “optimal” values for the tuning parameter α , where optimality is in terms of the right amount of data-specific downweighting needed to achieve robustness with minimal loss of efficiency. This is a valuable new contribution to the literature of DPD-based inference since, in contrast with the previous knowledge of α as trade-off between robustness and efficiency, our results lead to a small positive value of α as an optimum choice in practical applications. Furthermore, the estimation of random effects has been also considered, which is an important aspect in the study of LMMs. In particular, we provide a closed form to compute the random effects estimates based on the minimum DPD estimation. Large-scale numerical explorations, including an extensive simulation study, are provided to substantiate the theory developed and justify our claims of the superiority of the MDPDE in the proposed application domain of LMMs. Finally, our proposal is applied successfully (and robustly) to analyze two real-life data, one on orthodontic measures and another on foveal and extrafoveal vision acuity.

The rest of the paper is organized as follows. The MDPDE for non-homogeneous observations described in Ghosh and Basu (2013) is briefly presented in Sect. 2. In Sect. 3, we define the proposed estimator in case of linear mixed models, considering the estimation of fixed effect parameters and variances of random effects, and the prediction of random terms. The asymptotic and robustness properties of our procedures are considered in Sect. 4 together with the computational aspects and the case of balanced data. Section 5 reports the organization and results of the simulation study we conducted, comparing the performance of the MDPDE to the most recent methods, exploring also the case of contaminated data. Section 6 provides the

application of the proposed estimator to the real-data example on orthodontic measures. Concluding remarks are presented in Sect. 7. The proof of the main theorem is reported in Appendix A, while the Supplementary Material contains the derivation of equivariance properties, some additional theoretical and Monte Carlo results, an example which shows the robustness of predicted random terms, and the application to the real-data example on foveal and extrafoveal vision acuity.

2 The MDPDE for independent non-homogeneous observations

The density power divergence family was first introduced by Basu et al. (1998) as a measure of discrepancy between two probability density functions. The authors used this measure to robustly estimate the model parameters under the usual setup of independent and identically distributed data. The density power divergence measure $d_\alpha(g, f)$ between two probability densities g and f is defined, in terms of a single tuning parameter $\alpha \geq 0$, as

$$d_\alpha(g, f) = \int \left\{ f^{1+\alpha} - \left(1 + \frac{1}{\alpha}\right) f^\alpha g + \frac{1}{\alpha} g^{1+\alpha} \right\} \quad \text{if } \alpha > 0, \quad (1)$$

$$d_0(g, f) = \int g \ln \left(\frac{g}{f} \right) \quad \text{if } \alpha = 0, \quad (2)$$

where \ln denotes the natural logarithm. Basu et al. (1998) demonstrated that the tuning parameter α controls the trade-off between efficiency and robustness of the resulting estimator. With increasing α , the estimator acquires greater stability with a slight loss in efficiency. Since the divergence is not defined for $\alpha = 0$, $d_0(g, f)$ in Equation (2) represents the divergence obtained in the limit of (1) as $\alpha \rightarrow 0$, which corresponds to a version of the Kullback–Leibler divergence. On the other hand, $\alpha = 1$ generates the squared L_2 distance.

Let G be the true data generating distribution and g the corresponding density function. To model g , consider the parametric family of densities $\mathcal{F}_\theta = \{f_\theta : \theta \in \Theta \subseteq \mathbb{R}^p\}$. The minimizer of $d_\alpha(g, f_\theta)$ over $\theta \in \Theta$, whenever it exists, is the minimum DPD functional at the distribution point G . Note that the third term of the divergence $d_\alpha(g, f_\theta)$ is independent of θ ; hence, it can be discarded from the objective function as it has no role in the minimization process. Consider a sequence of independent and identically distributed (i.i.d) observations Y_1, \dots, Y_n from the true distribution G . Using the empirical distribution function G_n in place of G , the MDPDE of θ can be obtained by minimizing

$$\int f_\theta^{1+\alpha} - \left(1 + \frac{1}{\alpha}\right) \frac{1}{n} \sum_{i=1}^n f_\theta^\alpha(Y_i)$$

over $\theta \in \Theta$. In the above equation, the empirical distribution function is used to approximate its theoretical version (or, alternatively, the sample mean is used to

approximate the population mean). Note that, it is valid in case of continuous densities also. See Basu et al. (2011) for more details and examples.

Ghosh and Basu (2013) generalized the above concept of robust minimum DPD estimation to the more general case of independent non-homogeneous observations, i.e., they considered the case where the observed data Y_1, \dots, Y_n are independent but for each $i, Y_i \sim g_i$ where g_1, \dots, g_n are possibly different densities with respect to some common dominating measure. We model g_i by the family $\mathcal{F}_{i,\theta} = \{f_i(\cdot, \theta) : \theta \in \Theta\}$ for $i \in \{1, \dots, n\}$. While the distributions $f_i(\cdot, \theta)$ can be distinct, they share the same parameter vector θ . Ghosh and Basu (2013) proposed to minimize the average divergence between the data points and the model densities which leads to the minimization of the objective function

$$H_n(\theta) = \frac{1}{n} \sum_{i=1}^n \left[\int f_i(y, \theta)^{1+\alpha} dy - \left(1 + \frac{1}{\alpha}\right) f_i(Y_i, \theta)^\alpha \right] = \frac{1}{n} \sum_{i=1}^n H_i(Y_i, \theta), \quad (3)$$

where $H_i(Y_i, \theta)$ is the indicated term within the square brackets in the above equation. Differentiating the above expression with respect to θ , we get the estimating equations of the MDPDE for non-homogeneous observations. Note that the estimating equation is unbiased when each g_i belongs to the model family $\mathcal{F}_{i,\theta}$, respectively. When $\alpha \rightarrow 0$, the corresponding objective function reduces to $-\sum_{i=1}^n \ln(f_i(Y_i, \theta))/n$, which is the negative of the log-likelihood function. In Section SM-1 of the Supplementary Material, we report Assumptions (A1)–(A7) which are used to prove the asymptotic normality of the MDPDE (Ghosh and Basu 2013).

3 The MDPDE for linear mixed models

The general formulation of a LMM may be expressed as

$$Y = X\beta + ZU + \epsilon, \quad (4)$$

where $Y \in \mathbb{R}^d$ is the response vector, X and Z are known design matrices, β is the parameter vector for fixed effects, U is the vector of random effects, and ϵ is the random error vector. The vector U is assumed to be a random variable, in particular $U \sim N_q(\mathbf{0}, D)$, then $E(Y|U = u) = X\beta + Zu$. Notice that we will use U to indicate the random variable and u for its realization. Finally, assume that $\epsilon \sim N_d(\mathbf{0}, R)$ and ϵ and U are independent of each other, then

$$Y \sim N_d(X\beta, ZDZ^T + R). \quad (5)$$

The estimation of the matrices D and R involves a large amount of parameters; indeed, we need to assume some additional structure to the mixed model. According to the variance components model, the levels of any random effect are assumed to be independent with the same variance. Different effects are assumed independent with possibly different variances. Let the model have r random factors U_j with q_j levels, $j \in \{1, \dots, r\}$, with $q = \sum_{j=1}^r q_j$. As stated in Christensen (2011), a ‘‘natural generalization’’ of model 5 is to partition the vector $U = [U_1, \dots, U_r]$ and matrix

$\mathbf{Z} = [\mathbf{Z}_1 \dots \mathbf{Z}_r]$, such that \mathcal{U}_j are independent between each other, i.e., $\text{Cov}(\mathcal{U}_i, \mathcal{U}_j) = 0$ for $i \neq j$, and $\text{Cov}(\mathcal{U}_j) = \mathbf{D}_j = \sigma_j^2 \mathbf{I}_{q_j}$, $j = 1, \dots, r$. Finally, assume ϵ and \mathcal{U}_j independent for all j and $\mathbf{R} = \sigma_0^2 \mathbf{I}_d$, where \mathbf{I}_n is the $n \times n$ identity matrix.

Then, $\mathbf{Y} \sim N_d(\mathbf{X}\boldsymbol{\beta}, \mathbf{V})$ where

$$\mathbf{V} = \sigma_0^2 \mathbf{I}_d + \sum_{j=1}^r \mathbf{Z}_j \mathbf{Z}_j^\top \sigma_j^2 = \sigma_0^2 \left(\mathbf{I}_d + \sum_{j=1}^r \mathbf{Z}_j \mathbf{Z}_j^\top \gamma_j \right), \quad \text{with } \gamma_j = \frac{\sigma_j^2}{\sigma_0^2}.$$

In case of multiple measurements on each of a collection of observational units $\mathbf{Y}_1, \dots, \mathbf{Y}_n$, where n_i denotes the size of each group and $\sum_{i=1}^n n_i = N$, the LMM in equation (4) can be written for each observation as

$$\mathbf{Y}_i = \mathbf{X}_i \boldsymbol{\beta} + \sum_{j=1}^r \mathbf{Z}_{ij} \mathcal{U}_j + \epsilon_i, \quad i \in \{1, \dots, n\}, \tag{6}$$

where $\mathbf{Y}_i (n_i \times 1)$ is the response vector for group i , $\mathbf{X}_i (n_i \times k)$ and $\mathbf{Z}_{ij} (n_i \times q_j)$ are the model matrices, $\boldsymbol{\beta} (k \times 1)$ is the vector of unknown parameters for fixed effects, and $\epsilon_i (n_i \times 1)$ is the error term. Then, assuming $\epsilon_i \sim \sigma_0^2 \mathbf{I}_{n_i}$ and independent with respect to \mathcal{U}_j for all i, j ,

$$\mathbf{Y}_i \sim N_{n_i}(\mathbf{X}_i \boldsymbol{\beta}, \mathbf{V}_i), \quad i \in \{1, \dots, n\}.$$

where $\mathbf{V}_i = \sigma_0^2 (\mathbf{I}_{n_i} + \sum_{j=1}^r \mathbf{Z}_{ij} \mathbf{Z}_{ij}^\top \gamma_j)$.

Remark Note that the covariance structure of the linear mixed models considered here also includes the case of standard random intercept and random slope models. The real-data application presented in Sect. 6 provides an example. Indeed, Christensen (Christensen 2011) also considered this model for the discussion of all the variance components estimation methods. Also note that all the methods and results described in this paper can be routinely derived for any other (low-dimensional) parametric structure specified for the variance components in the LMM as per the need, but such situations (beyond the structure considered here) would rarely occur in practice.

3.1 Robust parameter estimation

In this setting, we can obtain the MDPDE for the p -dimensional parameter vector $\boldsymbol{\theta} = (\boldsymbol{\beta}^\top; \sigma_j^2, j \in \{0, \dots, r\})^\top$, with $p = k + r + 1$, by minimizing the objective function given in Equation (3) with $f_i \equiv N_{n_i}(\mathbf{X}_i \boldsymbol{\beta}, \mathbf{V}_i)$. Upon simplification, the objective function is given by

$$H_n(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n \left[\eta_{i\alpha} (\alpha + 1)^{-\frac{n_i}{2}} - \left(1 + \frac{1}{\alpha} \right) w_i(\boldsymbol{\theta}) \right] \tag{7}$$

where

$$\eta_{i\alpha} = (2\pi)^{-\frac{n_i\alpha}{2}} |\mathbf{V}_i|^{-\frac{\alpha}{2}}$$

and

$$w_i(\boldsymbol{\theta}) = \eta_{i\alpha} \exp \left\{ -\frac{\alpha}{2} (\mathbf{Y}_i - \mathbf{X}_i\boldsymbol{\beta})^\top \mathbf{V}_i^{-1} (\mathbf{Y}_i - \mathbf{X}_i\boldsymbol{\beta}) \right\}. \tag{8}$$

Differentiating the above equation with respect to $\boldsymbol{\beta}$, we get the corresponding estimating equation for the MDPDE of $\boldsymbol{\beta}$ as

$$\frac{\partial H_n}{\partial \boldsymbol{\beta}} = \frac{1}{n} \sum_{i=1}^n \left[- (1 + \alpha) w_i(\boldsymbol{\theta}) \mathbf{X}_i^\top \mathbf{V}_i^{-1} (\mathbf{Y}_i - \mathbf{X}_i\boldsymbol{\beta}) \right] = 0. \tag{9}$$

Let \mathbf{U}_{ij} denote the partial derivative of the matrix \mathbf{V}_i with respect to σ_j^2 . We have that $\mathbf{U}_{i0} = \mathbf{I}_{n_i}$, and $\mathbf{U}_{ij} = \mathbf{Z}_{ij} \mathbf{Z}_{ij}^\top$, $j \in \{1, \dots, r\}$. Then, the partial derivative of the objective function with respect to σ_j^2 , $j \in \{0, \dots, r\}$, leads to their MDPDE estimating equations as given by

$$\begin{aligned} \frac{\partial H_n}{\partial \sigma_j^2} = & \frac{1}{n} \sum_{i=1}^n \left\{ - \frac{\alpha \eta_{i\alpha} \text{Tr}(\mathbf{V}_i^{-1} \mathbf{U}_{ij})}{2(\alpha + 1)^{\frac{n_i}{2}}} + \frac{\alpha}{2} \left(1 + \frac{1}{\alpha} \right) w_i(\boldsymbol{\theta}) \right. \\ & \left. \times \left[\text{Tr}(\mathbf{V}_i^{-1} \mathbf{U}_{ij}) - (\mathbf{Y}_i - \mathbf{X}_i\boldsymbol{\beta})^\top \mathbf{V}_i^{-1} \mathbf{U}_{ij} \mathbf{V}_i^{-1} (\mathbf{Y}_i - \mathbf{X}_i\boldsymbol{\beta}) \right] \right\} = 0, \end{aligned} \tag{10}$$

where $\text{Tr}(\cdot)$ denotes the trace of the argument matrix. Note that for the case $\alpha = 0$, the estimating equations (9)–(10) correspond to the MLE score equations. Thus, the MDPDE at $\alpha = 0$ is nothing but the usual MLE under LMMs.

3.2 Robust estimation of the random effects

Finally, the estimates obtained by minimizing the objective function in (7) are used to predict the random realizations $\mathbf{u}_i | \mathbf{Y}_i$, $i = 1, \dots, n$. Consider the joint distribution of $(\mathbf{Y}_i, \mathbf{U}) \sim g_i(\mathbf{y}, \mathbf{u})$, and we have that $g_i(\mathbf{y}, \mathbf{u}) = g_i(\mathbf{y} | \mathbf{u}) g(\mathbf{u})$, where $g(\mathbf{u})$ is the true density of \mathbf{U} . Let $f_i(\mathbf{y}, \mathbf{u})$ be the parametric density to model $g_i(\mathbf{y}, \mathbf{u})$, then it can be expressed as $f_i(\mathbf{y}, \mathbf{u}) = f_i(\mathbf{y} | \mathbf{u}) f(\mathbf{u})$, which, by the plug-in principle, can be estimated using

$$f_i(\mathbf{y} | \mathbf{u}) = N_{n_i}(\mathbf{X}_i \hat{\boldsymbol{\beta}} + \mathbf{Z}_i \mathbf{u}_i, \hat{\sigma}_0^2 \mathbf{I}_{n_i}) \text{ and } f(\mathbf{u}) = N_q(\mathbf{0}, \hat{\mathbf{D}}).$$

Hence, we can estimate the random coefficients by minimizing the density power divergence measure between the densities $g_i(\mathbf{y}, \mathbf{u})$ and $f_i(\mathbf{y}, \mathbf{u})$ given by

$$H_n(\mathbf{u}_1, \dots, \mathbf{u}_n) = \frac{1}{n} \sum_{i=1}^n \kappa_\alpha \left[(\alpha + 1)^{-\frac{n_i+q}{2}} + \left(1 + \frac{1}{\alpha} \right) c_i(\mathbf{u}_i) \right],$$

where

$$c_i(\mathbf{u}_i) = \exp \left\{ -\frac{\alpha}{2\hat{\sigma}_0^2} (\mathbf{Y}_i - \mathbf{X}_i\hat{\boldsymbol{\beta}} - \mathbf{Z}_i\mathbf{u}_i)^\top (\mathbf{Y}_i - \mathbf{X}_i\hat{\boldsymbol{\beta}} - \mathbf{Z}_i\mathbf{u}_i) - \frac{\alpha}{2} \mathbf{u}_i^\top \hat{\mathbf{D}}^{-1} \mathbf{u}_i \right\}$$

and

$$\kappa_\alpha = \frac{1}{(2\pi)^{\frac{(n_i+q)\alpha}{2}} (\hat{\sigma}_0^2)^{\frac{n_i\alpha}{2}} |\hat{\mathbf{D}}|^{\frac{\alpha}{2}}}.$$

Differentiating the above equation with respect to $\mathbf{u}_i, i = 1, \dots, n$, we obtain the estimating equations for the MDPDE of \mathbf{u}_i as

$$\frac{\partial H_n}{\partial \mathbf{u}_i} = \frac{(1 + \alpha)\kappa_\alpha c_i(\mathbf{u}_i)}{n} \left[(\hat{\sigma}_0^2)^{-1} \mathbf{Z}_i^\top (\mathbf{Y}_i - \mathbf{X}_i\hat{\boldsymbol{\beta}} - \mathbf{Z}_i\mathbf{u}_i) - \hat{\mathbf{D}}^{-1} \mathbf{u}_i \right] = 0. \tag{11}$$

The computational aspects about how equations (9), (10) and (11) are solved numerically are treated in subsection 4.3.

Remark The proposed estimating equations given in equation (11) for the random effects \mathbf{u}_i can be simplified to

$$\left[(\hat{\sigma}_0^2)^{-1} \mathbf{Z}_i^\top (\mathbf{Y}_i - \mathbf{X}_i\hat{\boldsymbol{\beta}} - \mathbf{Z}_i\mathbf{u}_i) - \hat{\mathbf{D}}^{-1} \mathbf{u}_i \right] = 0$$

since κ_α does not depend on \mathbf{u}_i and $c_i(\mathbf{u}_i) > 0 \forall i$. Then, the estimating equation depends on α only through the estimates of the unknown parameters $\hat{\boldsymbol{\beta}}, \hat{\sigma}_0$ and $\hat{\mathbf{D}}$. For $\alpha \rightarrow 0$, the estimates of the unknown parameters correspond to the MLE, then the random effects predictions \hat{u}_i correspond to the standard predictions of random effects. Hence, they are the *Best Linear Unbiased Predictors* (BLUPs). Our formula differs from the standard formula since we provide estimates of the random effects for each i -th observation, while they are usually reported for each random effect.

4 Theoretical properties of the MDPDE under the Linear mixed models

4.1 Asymptotic efficiency

In this subsection, we prove that the theorem stated in Ghosh and Basu (2013), about the asymptotic behavior of the MDPDE for non-homogeneous observations, holds for the LMM setup. In particular, we present some conditions on the independent variables and the variance–covariance matrices that are used to derive the asymptotic distribution of the estimator in the LMM application.

We assume that the true densities $g_i, i = 1, \dots, n$, belong to the model family, i.e., $g_i = f_i(\cdot, \boldsymbol{\theta})$ for some value of $\boldsymbol{\theta} \in \Theta$. Consider the following assumptions.

(MM1) Define $X'_i = \left(\frac{\eta_{i\alpha} V_i^{-1}}{(1+\alpha)^{\frac{n_i}{2}+1}} \right)^{\frac{1}{2}} X_i$, for each i , and $X' = \text{Block-Diag} (X'_i : i \in \{1, \dots, n\})$.

Then, the X' matrix satisfies

$$\inf_n \left[\min \text{ eigenvalue of } \frac{X'^T X'}{n} \right] > 0, \tag{12}$$

and X_i and Z_i are full rank matrices for all i .

(MM2) The values of X_i 's are such that, for all j, k, l

$$\sup_{n>1} \max_{1 \leq i \leq n} |X_{ij}^T V_i^{-\frac{1}{2}}| = O(1), \quad \sup_{n>1} \max_{1 \leq i \leq n} |X_{ij}^T V_i^{-1} X_{ik}| = O(1), \tag{13}$$

$$\frac{1}{n} \sum_{i=1}^n |X_{ij}^T V_i^{-1} X_{ik} X_{il}^T V_i^{-\frac{1}{2}} \mathbf{1}| = O(1), \tag{14}$$

$$\frac{1}{n} \sum_{i=1}^n |X_{ij}^T V_i^{-\frac{1}{2}} |diag(V_i^{-\frac{1}{2}} X_{ik} X_{il}^T V_i^{-\frac{1}{2}}) \mathbf{1}| = O(1),$$

where $\mathbf{1}(n_i \times 1)$ is a vector of 1's.

(MM3) The matrices V_i and U_{ij} are such that, for all $j, k, l \in \{0, \dots, r\}$,

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n \text{Tr}(V_i^{-1} U_{ij}) = O(1), \quad \frac{1}{n} \sum_{i=1}^n \text{Tr}(V_i^{-1} U_{ij}) \text{Tr}(V_i^{-1} U_{ik}) = O(1), \\ \frac{1}{n} \sum_{i=1}^n \text{Tr}(V_i^{-1} U_{ik} V_i^{-1} U_{ij}) = O(1), \end{aligned} \tag{15}$$

$$\frac{1}{n} \sum_{i=1}^n \text{Tr}(V_i^{-1} U_{ij} V_i^{-1} U_{ik} V_i^{-1} U_{il}) = O(1),$$

$$\frac{1}{n} \sum_{i=1}^n \text{Tr}(V_i^{-1} U_{ij} V_i^{-1} U_{ik}) \text{Tr}(V_i^{-1} U_{il}) = O(1), \tag{16}$$

$$\frac{1}{n} \sum_{i=1}^n \text{Tr}(V_i^{-1} U_{ij}) \text{Tr}(V_i^{-1} U_{ik}) \text{Tr}(V_i^{-1} U_{il}) = O(1),$$

where the determinant $|V_i|$ is bounded away from both zero and infinity $\forall i$.

(MM4) Define $X_i^* = \left(\frac{\eta_{i\alpha}^2 V_i^{-1}}{(1+2\alpha)^{\frac{n_i}{2}+1}} \right)^{\frac{1}{2}} X_i$, for each i , and $X^* = \text{Block-Diag} (X_i^* : i \in \{1, \dots, n\})$.

Then, the X^* matrix satisfies

$$\max_{1 \leq i \leq n} \left[\frac{(X^{*T} X^*)^{-1} X_i^T V_i^{-1} X_i}{n} \right] = O(1). \tag{17}$$

Theorem 1 Consider the setup of the linear mixed model presented in Sect. 3. Assume that the true data generating density belongs to the model family and that the independent variables satisfy Assumptions (MM1)–(MM4) for a given (fixed) $\alpha \geq 0$. Then, we have the following results as $n \rightarrow \infty$ keeping n_i fixed for each i .

- (i) There exists a consistent sequence of roots $\hat{\theta}_n = (\hat{\beta}_n^\top, \hat{\sigma}_{n_j}^2, j \in \{0, \dots, r\})^\top$ to the minimum DPD estimating equations given in (9)–(10).
- (ii) The asymptotic distributions of $\hat{\beta}$ and $\hat{\sigma}_j^2$ are independent for all $j \in \{0, \dots, r\}$.
- (iii) The asymptotic distribution of $\Omega_n^{-\frac{1}{2}} \Psi_n \sqrt{n}(\hat{\theta}_n - \theta)$ is p -dimensional normal with mean zero and covariance matrix I_p . In particular, the asymptotic distribution of $(X^{*\top} X^*)^{-\frac{1}{2}} (X'^\top X')(\hat{\beta} - \beta)$ is a k -dimensional normal with mean zero and covariance matrix I_k , where X' and X^* are as defined in Assumptions (MM1) and (MM4), respectively.

The derivation of matrices Ω_n and Ψ_n and the proof of Theorem 1 are presented in Section A.1 of Appendix.

4.2 Influence function

To explore the robustness properties of the coefficient estimates in our treatment of linear mixed models, we derive the influence function of the MDPDEs. Denote the density power divergence functional $T_\alpha = (T_\alpha^\beta, T_\alpha^\Sigma)$ for the parameter vector $\theta^\top = (\beta^\top, \Sigma = (\sigma_0^2, \dots, \sigma_r^2))$. We continue with the notation of the previous subsections.

The influence function of the estimator T_α^β with contamination at the direction i_0 at the point t_{i_0} is computed to have the form

$$IF_{i_0}(t_{i_0}, T_\alpha^\beta, G_1, \dots, G_n) = (X'^\top X')^{-1} X_{i_0}^\top V_{i_0}^{-1} (t_{i_0} - X_{i_0} \beta) f_{i_0}(t_{i_0}; \theta)^\alpha, \tag{18}$$

and the corresponding influence function for the estimator T_α^Σ has the form

$$IF_{i_0}(t_{i_0}, T_\alpha^\Sigma, G_1, \dots, G_n) = \left[\sum_{i=1}^n \frac{\eta_{i\alpha}}{4(1 + \alpha)^{\frac{n_i}{2} + 2}} T(V_i^{-1} U_{ij}, V_i^{-1} U_{ik}) \right]^{-1} \tau_{i_0}, \tag{19}$$

where

$$\tau_i = \left(\frac{1}{2} f_i(t_i; \theta)^\alpha [\text{Tr}(V_i^{-1} U_{ij}) - (t_i - X_i \beta)^\top V_i^{-1} U_{ij} V_i^{-1} (t_i - X_i \beta)] - \frac{\eta_{i\alpha} \alpha \text{Tr}(V_i^{-1} U_{ij})}{2(1 + \alpha)^{\frac{n_i}{2} + 1}} \right).$$

Note that the influence functions in equations (18) and (19), seen as functions of the point t , are bounded for any $\alpha > 0$ since they are proportional to the functions ze^{-z^2}

and $\mathbf{z}^\top \mathbf{z} e^{-\mathbf{z}^\top \mathbf{z}}$, respectively, which are bounded for $\mathbf{z} \in \mathbb{R}^{n_i}$. For $\alpha = 0$, the influence functions for T_α^β and T_α^Σ are seen to be unbounded; indeed, this case corresponds to the non-robust maximum likelihood estimator. Hence, unlike the MLE, the minimum DPD estimators are B-robust, i.e., their associated influence functions are bounded, for $\alpha > 0$.

The influence function of the estimators T_α^β and T_α^Σ with contamination in all the n cases at the contamination points $\mathbf{t}_1, \dots, \mathbf{t}_n$, can be derived with similar computations and correspond to the sum of equations (18) and (19), respectively, for $i_0 = 1, \dots, n$. In this case also the influence functions are bounded for $\alpha > 0$ and unbounded for $\alpha = 0$.

Several summary measures of robustness based on the influence function for i.i.d. observations have been introduced in Hampel (1968, 1974). Following the same approaches, some influence function-based gross summary measures can be defined for the non-homogeneous case. For $\alpha > 0$, the gross-error sensitivity and the self-standardized sensitivity of the estimator T_α^β in the case of contamination only in the i_0 -th direction are given by

$$\begin{aligned} \gamma_{i_0}^u(T_\alpha^\beta, G_1, \dots, G_n) &= \sup_{\mathbf{t}} \{ |IF_{i_0}(\mathbf{t}, T_\alpha^\beta, G_1, \dots, G_n)| \} \\ &= \frac{\left[\lambda_{\max} \left((\mathbf{X}'^\top \mathbf{X}')^{-2} \mathbf{X}_{i_0}^\top \mathbf{V}_{i_0}^{-1} \mathbf{X}_{i_0} \right) \right]^{1/2}}{\sqrt{\alpha} (2\pi)^{\frac{n_{i_0}\alpha}{2}} |\mathbf{V}_{i_0}|^{\frac{\alpha}{2}} e^{1/2}} \end{aligned} \tag{20}$$

and

$$\begin{aligned} \gamma_{i_0}^s(T_\alpha^\beta, G_1, \dots, G_n) &= \sup_{\mathbf{t}} \{ IF_{i_0}(\mathbf{t}, T_\alpha^\beta, G_1, \dots, G_n)^\top (\Psi_n^{-1} \mathbf{\Omega}_n \Psi_n^{-1})^{-1} IF_{i_0}(\mathbf{t}, T_\alpha^\beta, G_1, \dots, G_n) \}^{\frac{1}{2}} \\ &= \frac{\left[\lambda_{\max} \left((\mathbf{X}^{*\top} \mathbf{X}^*)^{-1} \mathbf{X}_{i_0}^\top \mathbf{V}_{i_0}^{-1} \mathbf{X}_{i_0} \right) \right]^{1/2}}{n \sqrt{\alpha} (2\pi)^{\frac{n_{i_0}\alpha}{2}} |\mathbf{V}_{i_0}|^{\frac{\alpha}{2}} e^{1/2}}, \end{aligned} \tag{21}$$

where $\lambda_{\max}(\mathbf{A})$ indicates the largest eigenvalue of the matrix \mathbf{A} , while they are equal to ∞ if $\alpha = 0$. Details of the computations are provided in Section SM-2 of the Supplementary Materials. The sensitivity measures for T_α^Σ have no compact form and they are not reported.

4.3 Computational aspects

The estimating equations (9)–(10) previously introduced can be solved numerically in order to obtain the estimates of $\boldsymbol{\theta} = (\boldsymbol{\beta}^\top; \sigma_j^2, j \in \{0, \dots, r\})^\top$.

Note that denoting $w_j = w_j(\boldsymbol{\beta}, \sigma_j^2)$, $j = 0, \dots, r$, as defined in equation (8), the estimating equation for $\boldsymbol{\beta}$ given in equation (9) corresponds to

$$-\sum_{i=1}^n w_i X_i^T V_i^{-1} Y_i + \sum_{i=1}^n w_i X_i^T V_i^{-1} X_i \beta = 0.$$

Solving for β , we get

$$\beta = \left(\sum_{i=1}^n w_i X_i^T V_i^{-1} X_i \right)^{-1} \left(\sum_{i=1}^n w_i X_i^T V_i^{-1} Y_i \right),$$

so that, in an iterative fixed point algorithm, the successive iterates have the relation

$$\beta^{(k+1)} = \left(\sum_{i=1}^n w_i(\beta^{(k)}, \sigma_j^{(k)}) X_i^T V_i(\sigma_j^{(k)})^{-1} X_i \right)^{-1} \left(\sum_{i=1}^n w_i(\beta^{(k)}, \sigma_j^{(k)}) X_i^T V_i(\sigma_j^{(k)})^{-1} Y_i \right).$$

The estimating equation (10) for the variance components cannot be written in a similar closed form, then it is solved numerically by a quasi-Newton method which allows box constraints, that is, each variable can be given a lower and/or upper bound.

Finally, the random coefficients can be predicted through the obtained estimates. Notice that the quantity κ_α can be removed to solve equation (11) since it does not depend on u_i , as well as $c_i(u_i)$ since $c_i(u_i) > 0$ for all i . Hence, equation (11) is simplified to

$$(\hat{\sigma}_0^2)^{-1} Z_i^T (Y_i - X_i \hat{\beta}) - \left((\hat{\sigma}_0^2)^{-1} Z_i^T Z_i + \hat{D}^{-1} \right) u_i = 0.$$

Then,

$$\hat{u}_i = \left((\hat{\sigma}_0^2)^{-1} Z_i^T Z_i + \hat{D}^{-1} \right)^{-1} (\hat{\sigma}_0^2)^{-1} Z_i^T (Y_i - X_i \hat{\beta}). \tag{22}$$

This formulation provides a closed form to predict the realizations $u_i, i = 1, \dots, n$. It is worth noting that even if these estimates may appear not robust, in the sense that an outliers (Y_i, X_i) can affect the predicted u_i , they are based on the robust estimates $\hat{\beta}, \hat{\sigma}_0^2$ and \hat{D} .

4.4 An example: the balanced data case

Consider the model defined by Equation (6). Here, we study the simplest case in which $n_i = p$, for all $i \in \{1, \dots, n\}$, and the associated random effects covariates (Z_{ij}) are also the same for all i . In this case, the covariance matrix of Y_i is the same for all i and is denoted by V having the form

$$V = V_i = \sigma_0^2 \left(I_p + \sum_{j=1}^r U_j \gamma_j \right),$$

where $\gamma_j = \sigma_j^2 / \sigma_0^2$ and $U_j = U_{ij}$ as it is independent of i .

The influence function of the functional T_α^β with contamination in the direction i_0 , given in Equation (18), can be written as

$$IF_{i_0}(t_{i_0}, T_\alpha^\beta, G_1, \dots, G_n) = (1 + \alpha)^{\frac{p}{2}+1} \left(\sum_{i=1}^n X_i^T V^{-1} X_i \right)^{-1} X_{i_0}^T V^{-1} (t_{i_0} - X_{i_0} \beta) w_{i_0}.$$

Using this expression, the gross-error sensitivity for the functional T_α^β is given by

$$\gamma_{i_0}^u(T_\alpha^\beta) = \frac{(1 + \alpha)^{\frac{p}{2}+1}}{\sqrt{\alpha}} \left[\lambda_{\max} \left(\left[\sum_{i=1}^n X_i^T V^{-1} X_i \right]^{-2} X_{i_0}^T V^{-1} X_{i_0} \right) \right]^{1/2} e^{-1/2}. \tag{23}$$

Similarly, the self-standardized sensitivity of the functional T_α^β can be written as

$$\gamma_{i_0}^s(T_\alpha^\beta) = \frac{(1 + \alpha)^{\frac{p+2}{4}}}{n\sqrt{\alpha}} \left[\lambda_{\max} \left(\left[\sum_{i=1}^n X_i^T V^{-1} X_i \right]^{-1} X_{i_0}^T V_{i_0}^{-1} X_{i_0} \right) \right]^{1/2} (2e)^{-1/2}. \tag{24}$$

The simpler form of the sensitivity measures allows us to assess their performance with respect to the tuning parameter α . Indeed, the function $\frac{(1+\alpha)^{\frac{p}{2}+1}}{\sqrt{\alpha}}$ in the gross-error sensitivity (23) has a minimum for the value $\alpha^* = \frac{1}{p+1}$, suggesting that this value of the parameter α gives the most robust estimator. Similarly, the function $\frac{(1+\alpha)^{\frac{p+2}{4}}}{\sqrt{\alpha}}$ in the self-standardized sensitivity (24) has a minimum for the value $\bar{\alpha} = \frac{2}{p}$. The existence of such minimum values for α is in contrast with the previously held knowledge about this parameter. It was introduced as a trade-off between efficiency and robustness, instead here we show that, passing some threshold, with increasing α , we lose both efficiency and robustness. Finally, the proposed optimal values α^* and $\bar{\alpha}$ depend only on the dimension of observations and constitute valuable choices in practical situations.

In the following, we present a simple example for which we will compute the theoretical quantities introduced above. This example in linear mixed models has been chosen for its similarity to the case of longitudinal data; it is often also named as LMM with random intercept and random slope.

Notice that the MDPDE satisfies the equivariance properties. In particular, the regression equivariance allows us to assume, without loss of generality, any suitable value for the parameter β while proving the asymptotic properties with the following example or for the Monte Carlo studies. Since the MDPDE is also scale and affine equivariant, such estimators do not depend on the choice of the coordinate system for the variables x and on the measurement unit of y . The derivation of these properties is reported in Section SM-3 of the Supplementary Material. Furthermore, Section SM-4 reports an illustrative example which shows the robustness achieved by the random effects predictions computed according to equation (22).

We consider $n = 50$ different subjects (groups), and for each of them, we have $p = 10$ measurements taken with respect to the factor u_{i2} , $i = 1, \dots, n$, with two levels, modeled here as a random effect. The \mathbf{X} 's model matrices are simulated from a standard normal. In particular, the model is described by

$$\mathbf{Y}_i = \beta_0 + \beta_1 \mathbf{X}_i + \mathbf{U}_1 + \mathbf{U}_2 \mathbf{Z}_{i2} + \epsilon_i,$$

where $\mathbf{U}_1 \sim N_p(0, \sigma_1^2 \mathbf{I}_p)$, $\mathbf{U}_2 \sim N_2(0, \sigma_2^2 \mathbf{I}_2)$ and $\epsilon_i \sim N_p(0, \sigma_0^2 \mathbf{I}_p)$, and they are independent. Hence, for this model, $\boldsymbol{\theta} = (\beta_0, \beta_1, \sigma_1^2, \sigma_2^2, \sigma_0^2)$, and we take $\boldsymbol{\theta} = (1, 2, 0.25, 0.5, 0.25)$ as the true values of the parameters.

Using the given values, we compute the variance–covariance matrices \mathbf{V}_i and the matrices $\boldsymbol{\Psi}_n$ and $\boldsymbol{\Omega}_n$. First, we will look at the *Asymptotic Relative Efficiency* (ARE) of the minimum density power divergence estimators with respect to the fully efficient maximum likelihood estimator. Figure 1 shows the asymptotic relative efficiencies of the estimators of β_1 and σ_2^2 for $\alpha \in [0, 0.6]$. It is easy to see that there is a loss of efficiency which increases with α . However, for small positive values of α , the estimator retains reasonable efficiency. The ARE of the estimators of the other parameters are similar to those displayed here, and are given in Section SM–5 of the Supplementary Materials.

On the other hand, to study the robustness properties, Fig. 2 shows the influence functions of $T_\alpha^{\beta_0}$ and $T_\alpha^{\sigma_1^2}$, with respect to $\alpha = 0, 0.05, \alpha^*, \bar{\alpha}$ where $\alpha^* = 1/(p+1) = 1/11$ and $\bar{\alpha} = 2/p = 0.2$. Here, we have plotted $IF(\mathbf{t}_1, \dots, \mathbf{t}_n, T_\alpha, G_1, \dots, G_n)$, the influence function of the estimator T_α , computed with respect to constant vectors $\mathbf{t} = t(1, \dots, 1)^\top$ for varying $t \in \mathbb{R}$. Note that except for the case $\alpha = 0$, we can easily see that the influence function is bounded as may also be noted from equations (18) and (19); thus, the estimator will be robust with respect to outliers. The influence function for the estimators of other parameters behaves similarly; these plots are available in Section SM–5 of the Supplementary Materials.

Finally, Fig. 3 shows the gross-error sensitivity and the self-standardized sensitivity of the functional T_α^β . Here, we have considered a particular direction $i_0 \in \{1, \dots, n\}$. Note that in the present case of balanced data, the choice of i_0 does not change the behavior of the sensitivity measures with respect to α .

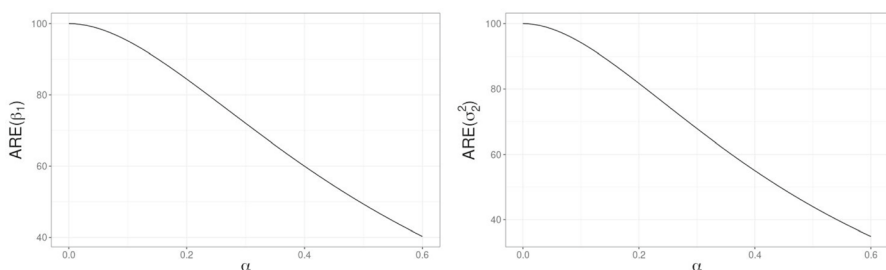
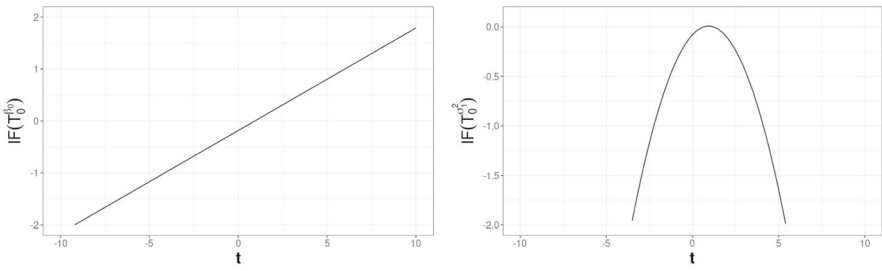
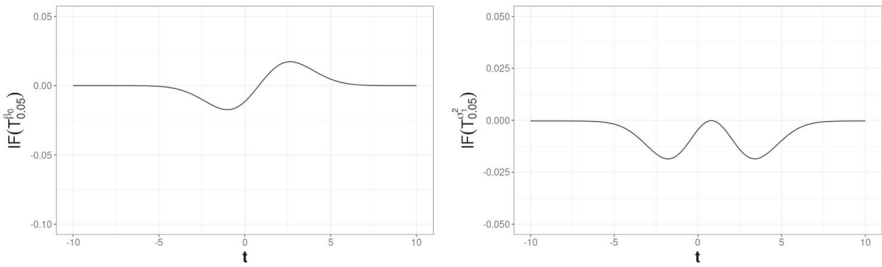


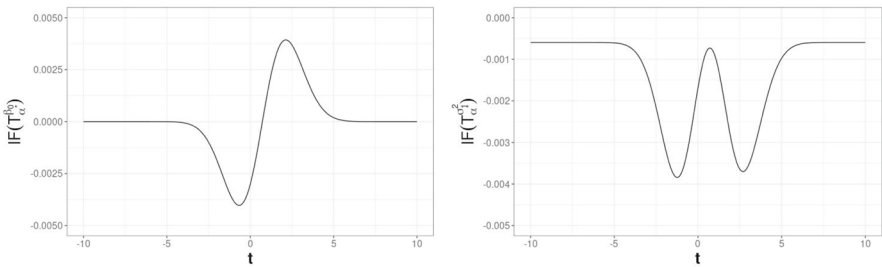
Fig. 1 Asymptotic relative efficiency with respect to α for β_1 and σ_2^2 , respectively



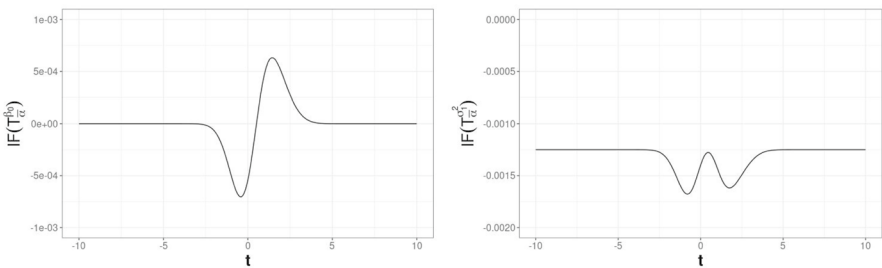
(a) $\alpha = 0$



(b) $\alpha = 0.05$



(c) $\alpha = \alpha^* = \frac{1}{p+1} = \frac{1}{11}$



(d) $\alpha = \bar{\alpha} = \frac{2}{p} = 0.2$

Fig. 2 Influence function for the functionals $T_\alpha^{\beta_0}$ (left panel) and $T_\alpha^{\sigma_1^2}$ (right panel), for different values of the tuning parameter α

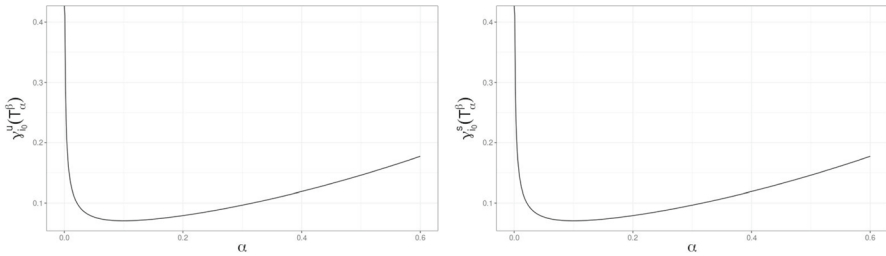


Fig. 3 Gross-error sensitivity (left panel) and self-standardized sensitivity (right panel) of the functional T_α^β with respect to $i_0 = 10$

5 Monte Carlo simulations

5.1 Model setting

We consider a simulation setting as introduced in Agostinelli and Yohai (2016) and reported here in order to facilitate the comparison of the considered estimators.

Consider an LMM for a two-way cross classification with interaction, where the model is given by

$$y_{fgh} = \mathbf{x}_{fgh}^\top \boldsymbol{\beta}_0 + a_f + b_g + c_{fg} + e_{fgh},$$

where $f = 1, \dots, F, g = 1, \dots, G,$ and $h = 1, \dots, H.$ Here, we set $F = 2, G = 2$ and $H = 3$ getting $p = F \times G \times H = 12.$ Also \mathbf{x}_{fgh} is a $k \times 1$ vector where the last $k - 1$ components are from a standard multivariate normal and the first component is identically equal to 1, and $\boldsymbol{\beta}_0 = (0, 2, 2, 2, 2, 2)^\top$ is a $k \times 1$ vector of the fixed parameters with $k = 6.$ The random variables a_f, b_g and c_{fg} are the random effects which are normally distributed with variances $\sigma_a^2, \sigma_b^2,$ and $\sigma_c^2.$ Arranging the y_{fgh} in lexicon order (ordered by h within g within f), we obtain the vector \mathbf{y} of dimension $p,$ and in the similar way, the $p \times k$ matrix \mathbf{x} obtained arranging $\mathbf{x}_{fgh}.$ Similarly, we set $\mathbf{a} = (a_1, \dots, a_F)^\top, \mathbf{b} = (b_1, \dots, b_G)^\top$ and $\mathbf{c} = (c_{11}, \dots, c_{FG})^\top,$ that is, $\mathbf{a} \sim N_F(\mathbf{0}, \sigma_a^2 \mathbf{I}_F)$ and similarly for \mathbf{b} and $\mathbf{c},$ while $\mathbf{e} = (e_{111}, \dots, e_{FGH})^\top \sim N_p(\mathbf{0}, \sigma_e^2 \mathbf{I}_p).$ Hence, \mathbf{y} is a p multivariate normal with mean $\boldsymbol{\mu} = \mathbf{x}\boldsymbol{\beta}$ and variance matrix $\boldsymbol{\Sigma}_0 = \boldsymbol{\Sigma}(\eta_0, \boldsymbol{\gamma}_0) = \eta_0(\mathbf{V}_0 + \sum_{j=0}^J \gamma_j \mathbf{V}_j),$ where $\mathbf{V}_0 = \mathbf{I}_p, \mathbf{V}_1 = \mathbf{I}_F \otimes \mathbf{J}_G \otimes \mathbf{J}_H, \mathbf{V}_2 = \mathbf{J}_F \otimes \mathbf{I}_G \otimes \mathbf{J}_H,$ and $\mathbf{V}_3 = \mathbf{I}_F \otimes \mathbf{I}_G \otimes \mathbf{J}_H;$ \otimes is the Kronecker product and \mathbf{J}_k is a $k \times k$ matrix of ones. We took $\sigma_a^2 = \sigma_b^2 = 1/16$ and $\sigma_c^2 = 1/8.$ Then, $\boldsymbol{\gamma}_0 = (\gamma_{01}, \gamma_{02}, \gamma_{03})^\top = (\sigma_a^2/\sigma_e^2, \sigma_b^2/\sigma_e^2, \sigma_c^2/\sigma_e^2)^\top = (1/4, 1/4, 1/2)^\top$ and $\eta_0 = \sigma_e^2 = 1/4.$

We consider a sample of size $n = 100$ and four levels of contamination $\varepsilon = 0, 5, 10$ and $15\%.$ Hence, $n \times \varepsilon$ observations are contaminated according the following contamination scenarios. Let \mathbf{y}_0 and \mathbf{x}_0 indicate the response vector and the fixed effect model matrix for the contaminated observations.

- Complete contamination: $n \times \varepsilon$ elements of the vector \mathbf{y} are replaced by observations from $\mathbf{y}_0 \sim N_p(\mathbf{x}_0 \boldsymbol{\beta}_0 + \boldsymbol{\omega}_0, \boldsymbol{\Sigma}).$ The matrix \mathbf{x}_0 is such that the

first column is identically equal to 1, while the last $k - 1$ columns are from $N_{p \times (k-1)}(\boldsymbol{\phi}_0, 0.005^2 \mathbf{I}_{p \times (k-1)})$ where $\boldsymbol{\phi}_0$ indicates a p -vector of constants all equal to ϕ_0 with $\phi_0 = 1, 20$ in the case of low leverage outliers (lev1) or for large leverage outliers (lev20), respectively. $\boldsymbol{\omega}_0$ is a p -vector of constants all equal to ω_0 with $\omega_0 = 0, 1, \dots, 30$.

- Contamination only on the response \mathbf{y} : $\mathbf{y}_0 \sim N_p(\mathbf{x}\boldsymbol{\beta}_0 + \boldsymbol{\omega}_0, \boldsymbol{\Sigma})$ with ω_0 as above.
- Contamination only on the \mathbf{x} : $n \times \varepsilon$ rows of \mathbf{x} are replaced by \mathbf{x}_0 , as defined in the first scenario, with $\phi_0 = 0, 1, 3, 5, 10, 15, 20, 25, 30$. The corresponding elements of \mathbf{y} are replaced by observations sampled from $\mathbf{y}_0 \sim N_p(\mathbf{x}_0\boldsymbol{\beta}_0, \boldsymbol{\Sigma})$.

For each contamination scenario, we compute the CVFS-estimator described in Copt and Victoria-Feser (2006) with Rocke ρ function and with asymptotic rejection probability set to 0.01 as implemented in the R R Core Team (2019) package `robustvarComp` (Agostinelli and Yohai (2019)), the SMDM estimator introduced by Koller (2013) as implemented in the R package `robustlmm` (Koller (2016)), the composite τ -estimator proposed by Agostinelli and Yohai (2016) and available in the R package `robustvarComp`, and our proposed MDPDE with different choices of α . In particular, $\alpha \in \{0, 0.01, 0.1, 0.2, \dots, 1\}$; note that $\alpha^* = 1/(p + 1) = 1/13$ and $\bar{\alpha} = 2/p = 1/6$. For each case, we run 500 Monte Carlo replications.

5.2 Performance Measures

Let (\mathbf{y}, \mathbf{x}) be an observation independent of the sample $(\mathbf{y}_1, \mathbf{x}_1), \dots, (\mathbf{y}_n, \mathbf{x}_n)$ used to compute $\hat{\boldsymbol{\beta}}$ and let $\hat{\mathbf{y}} = \mathbf{x}\hat{\boldsymbol{\beta}}$ be the predicted value of \mathbf{y} using \mathbf{x} . Then, the squared Mahalanobis distance between $\hat{\mathbf{y}}$ and \mathbf{y} using the matrix $\boldsymbol{\Sigma}_0$ is

$$\begin{aligned}
 m(\hat{\mathbf{y}}, \mathbf{y}, \boldsymbol{\Sigma}_0) &= (\hat{\mathbf{y}} - \mathbf{y})^\top \boldsymbol{\Sigma}_0^{-1} (\hat{\mathbf{y}} - \mathbf{y}) \\
 &= (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}_0)^\top \mathbf{x}^\top \boldsymbol{\Sigma}_0^{-1} \mathbf{x} (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}_0) + (\mathbf{y} - \mathbf{x}\boldsymbol{\beta}_0)^\top \boldsymbol{\Sigma}_0^{-1} (\mathbf{y} - \mathbf{x}\boldsymbol{\beta}_0).
 \end{aligned}$$

Since $\mathbf{y} - \mathbf{x}\boldsymbol{\beta}_0$ is independent of \mathbf{x} and has covariance matrix $\boldsymbol{\Sigma}_0$, putting $\mathbf{A} = \mathbb{E}(\mathbf{x}^\top \boldsymbol{\Sigma}_0^{-1} \mathbf{x})$, we have

$$\begin{aligned}
 \mathbb{E}[m(\hat{\mathbf{y}}, \mathbf{y}, \boldsymbol{\Sigma}_0)] &= \mathbb{E}\left[(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}_0)^\top \mathbf{A} (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}_0)\right] + \text{trace}\left[\boldsymbol{\Sigma}_0^{-1} (\mathbf{y} - \mathbf{x}\boldsymbol{\beta}_0)(\mathbf{y} - \mathbf{x}\boldsymbol{\beta}_0)^\top\right] \\
 &= \mathbb{E}\left[(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}_0)^\top \mathbf{A} (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}_0)\right] + p.
 \end{aligned}$$

Then, to evaluate an estimator $\hat{\boldsymbol{\beta}}$ of $\boldsymbol{\beta}$ by its prediction performance, we can use

$$\mathbb{E}\left[m(\hat{\boldsymbol{\beta}}, \boldsymbol{\beta}_0, \mathbf{A})\right] = \mathbb{E}\left[(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}_0)^\top \mathbf{A} (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}_0)\right].$$

Let N be the number of replications in the simulation study, and let $\hat{\boldsymbol{\beta}}_j, 1 \leq j \leq N$ be the value of $\hat{\boldsymbol{\beta}}$ at the j -th replication, then we can estimate $\mathbb{E}\left[m(\hat{\boldsymbol{\beta}}, \boldsymbol{\beta}_0, \mathbf{A})\right]$ by the mean square Mahalanobis distance as

$$\text{MSMD} = \frac{1}{N} \sum_{j=1}^N m(\hat{\beta}_j, \beta_0, \mathbf{A}).$$

It is easy to prove that as in this case, \mathbf{x} is a $p \times k$ matrix where the cells are independent $N(0, 1)$ random variables, then $\mathbf{A} = \text{trace}(\Sigma_0^{-1})\mathbf{I}_k$.

Given two p -dimensional covariance matrices Σ_1 and Σ_0 , one way to measure how close Σ_1 and Σ_0 are is through the use of the Kullback–Leibler divergence between two multivariate normal distributions with the same mean and covariance matrices equal to Σ_1 and Σ_0 , given by

$$\text{KLD}(\Sigma_1, \Sigma_0) = \text{trace}(\Sigma_1 \Sigma_0^{-1}) - \log(\det(\Sigma_1 \Sigma_0^{-1})) - p.$$

Since (η_0, γ_0) determines $\Sigma_0 = \Sigma(\eta_0, \gamma_0)$, the covariance matrix of \mathbf{y} given \mathbf{x} for the particular LMM considered in our simulation (as described in Sect. 5.1), one way to measure the performance of an estimator $(\hat{\eta}, \hat{\gamma})$ of (η_0, γ_0) is by $\mathbb{E}[\text{KLD}(\Sigma(\hat{\eta}, \hat{\gamma}), \Sigma_0)]$. Let $(\hat{\eta}_j, \hat{\gamma}_j)$, $1 \leq j \leq N$, be the value of $(\hat{\eta}, \hat{\gamma})$ at the j -th replication, then we can estimate $\mathbb{E}[\text{KLD}(\Sigma(\hat{\eta}, \hat{\gamma}), \Sigma_0)]$ by the mean Kullback–Leibler divergence

$$\text{MKLD} = \frac{1}{N} \sum_{j=1}^N \text{KLD}(\Sigma(\hat{\eta}_j, \hat{\gamma}_j), \Sigma_0).$$

5.3 Results

We begin with the performance of the estimators in the absence of contamination. Table 1 shows the relative efficiency of the CVFS-estimator, the SMDM-estimator, the Composite τ -estimator and the MDPDE for different values of α with respect to maximum likelihood. The efficiency of the estimators of β has been measured by the

Table 1 Relative efficiency for the SMDM-estimator, CVFS-estimator, Composite τ -estimator and MDPDE for different values of α with respect to the maximum likelihood computed by the MSMD for the fixed terms β and by the MKLD for the random terms

Method	(α)	MSMD EFF.	MKLD EFF.
SMDM	–	0.956	0.147
CVFS	–	0.706	0.453
Composite τ	–	0.807	0.833
MDPDE	0.01	0.999	0.996
	α^*	0.960	0.945
	0.1	0.937	0.915
	$\tilde{\alpha}$	0.853	0.814
	0.2	0.805	0.760
	0.3	0.658	0.603
	0.4	0.519	0.470
	0.5	0.400	0.361
	0.6	0.302	0.273

MSMD ratio, while the MKLD ratio was used for the efficiency of an estimator of (η, γ) .

The MDPDEs exhibit a high relative efficiency, even greater than the competitor estimators, for small values of α , while the efficiency decreases with increasing α . Note that the MDPDEs are far more successful in retaining the efficiency of the estimators of the random components. For very small values of α , the MDPDEs dominate either competitor (at least up to $\alpha = \alpha^*$ for SMDM, and at least up to $\alpha = 0.2$ for CVFS) in terms of both (MSMD and MKLD) efficiency measures. As the value of α increases, the MSMD efficiency of the MDPDE eventually lags behind its competitors, but in terms of MKLD efficiency, it beats the competitors at least up to $\alpha = 0.4$, except for the Composite τ -estimator that has a higher efficiency than the MDPDE for $\alpha > \bar{\alpha}$. On the whole, it is clear that under pure data, a properly chosen member of the MDPDE class can perform competitively, if not better, compared to the SMDM, the CVFS and the Composite τ -estimators.

Now, we consider the contamination settings. At first, Tables 2 and 3 report the maximum values of MSMD and MKLD over the values of ω_0 and leverage ϕ_0 considered for the complete contamination and separated contamination, respectively, of the MDPDE for different values of α compared to the CVFS-, the SMDM- and Composite τ -estimators.

Small values of α , as expected, provide much higher maximum values with respect to the other estimators in Table 2. However, for slightly larger values of α , the MDPDEs are extremely competitive with the existing estimators. It may be easily observed that in case of complete contamination, the MDPDE at $\bar{\alpha}$ clearly beats the competitors (CVFS, SMDM and Composite τ) over both performance measures at both leverage values (except in case of complete contamination at MSMD, lev20, where its performance measure is equal to that of CVFS). In this example, the MDPDE at $\alpha = 0.2$ fares even better. In case of contamination on the y , the

Table 2 Complete contamination. Maximum values of MSMD and MKLD for the CVFS-, SMDM-, Composite τ -estimators and for the MDPDE at different values of α under 10% of outlier contamination

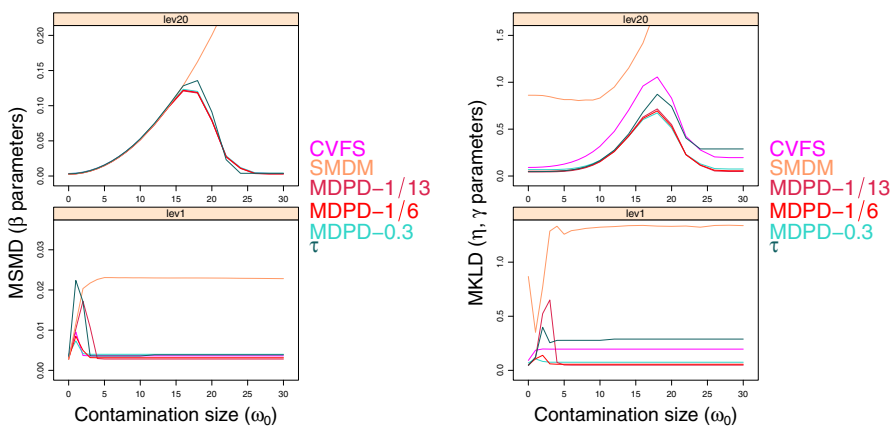
Method	(α)	MSMD		MKLD	
		lev1	lev20	lev1	lev20
CVFS	–	0.010	0.122	0.197	1.057
SMDM	–	0.023	0.450	1.341	9.125
Composite τ	–	0.022	0.136	0.398	0.872
MDPDE	0	9.007	9.005	3.508e22	1.605e25
	0.01	1.114	0.120	106.185	0.732
	$\alpha^*(\frac{1}{13})$	0.017	0.121	0.650	0.716
	0.1	0.012	0.121	0.387	0.710
	$\bar{\alpha}(\frac{1}{6})$	0.008	0.122	0.139	0.695
	0.2	0.008	0.122	0.105	0.688
	0.3	0.007	0.123	0.106	0.673
	0.4	0.007	0.125	0.116	0.665
	0.5	0.008	0.127	0.137	0.662
	0.6	0.010	0.130	0.170	0.666

Table 3 Contamination on x and y separately. Maximum values of MSMD and MKLD for the CVFS-, SMDM-, Composite τ -estimators and for MDPDE considering different values of α under 10% of outlier contamination

Method	(α)	MSMD		MKLD	
		x	y	x	y
CVFS	–	0.004	0.008	0.092	0.210
SMDM	–	0.003	0.021	0.820	1.667
Composite τ	–	0.004	0.010	0.042	0.305
MDPDE	0	0.003	8.998	1.012	203.248
	0.01	0.003	1.134	0.035	109.762
	$\alpha^*(\frac{1}{13})$	0.003	0.016	0.036	0.609
	0.1	0.003	0.011	0.037	0.389
	$\bar{\alpha}(\frac{1}{6})$	0.003	0.008	0.042	0.134
	0.2	0.003	0.007	0.046	0.109
	0.3	0.004	0.007	0.060	0.108
	0.4	0.005	0.007	0.081	0.120
	0.5	0.007	0.008	0.108	0.145
	0.6	0.009	0.009	0.145	0.185

estimators show similar performance, while, when only the x is contaminated, the MSMD values are almost equal for all the considered estimators suggesting that the estimation of the fixed effects parameter β is not affected. However, as in the previous case, the MDPDE for values of α close to zero show larger MKLD values, but slightly increasing α the MDPDE outperforms the competitors.

Figures 4a and 4b display the MSMD and MKLD as function of ω_0 , comparing the CVFS-, SMDM- and Composite τ -estimators with the MDPDEs for three chosen values of α , under 10% of outlier contamination. In particular, we choose



(a) MSMD performance of the estimators of β (b) MKLD performance of the estimators of (η, γ)

Fig. 4 Complete contamination. Performance of the MDPD-estimators of β and (η, γ) for $\alpha = \frac{1}{13}, \frac{1}{6}, 0.3$, compared to the CVFS-, SMDM- and Composite τ -estimators, under 10% outlier contamination

α^* and $\bar{\alpha}$ since they are the values suggested by theory, and $\alpha = 0.3$ since it shows the lowest (or very close to the lowest) maximum values of MSMD and MKLD. We can see that most of the MDPDEs outperform the CVFS- and SMDM-estimators, especially in case of leverage 20 (lev20), where the SMDM-estimator shows an unbounded behavior. On the other hand, in the case of leverage 1 (lev1), even if the CVFS-estimator presents lower maximum value of MSMD and MKLD for very small values of ω_0 , the MDPDEs show a better performance when ω_0 increases. In fact, the MDPDE at $\alpha = 0.3$ is competitive or better than CVFS at all values of ω_0 . Finally, Fig. 5a and 5b shows the MSMD and MKLD as function of ω_0 if the response y is contaminated and as function of ϕ_0 in case of contamination on the x . The CVFS-, SMDM- and Composite τ -estimators are compared with the MDPDEs for three chosen values of α , under 10% of outlier contamination. These results confirm that the MDPDEs, at least for one value of α , outperform the competitor estimators. Complete results of the simulation study, for all the considered estimators and levels of contamination, are reported in Section SM-6 of the Supplementary Materials.

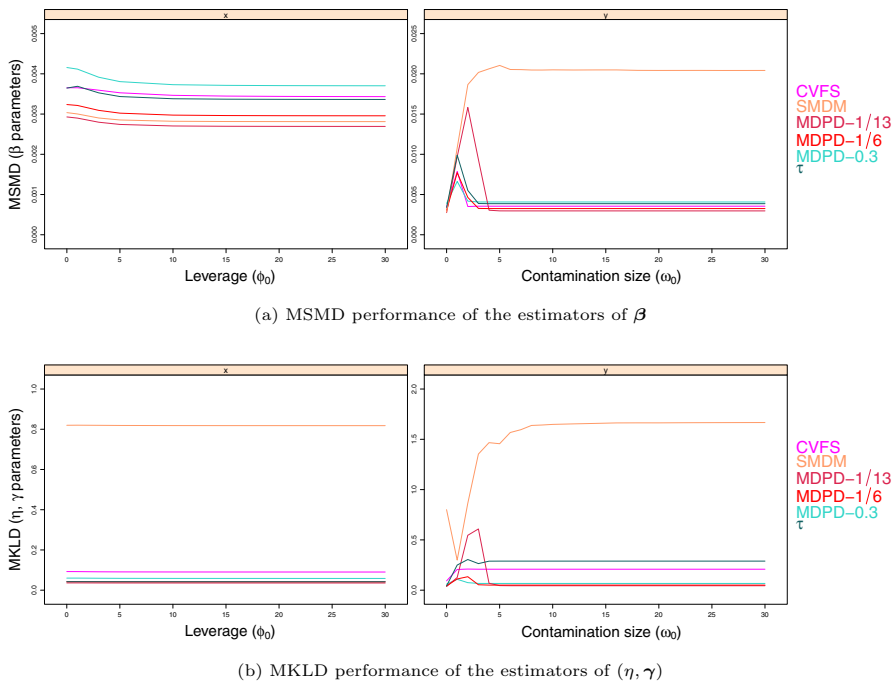


Fig. 5 Contamination on x and y . Performance of the MDPD-estimators of β and (η, γ) for $\alpha = \frac{1}{13}, \frac{1}{6}, 0.3$, compared to the CVFS-, SMDM- and Composite τ -estimators, under 10% outlier contamination

6 Real-data example: Orthodontic Distance Growth

Let us now present an application of the proposed estimation method to a real-data example on orthodontic measures, while Section SM–8 of the Supplementary Materials reports the analysis of the real-life data on foveal and extrafoveal vision acuity (and crowding) studying their interrelationships with one’s reading performances. We compare the estimates obtained by the minimum DPD method with those obtained using the classical (non-robust) restricted MLE, computed using the `lme` function in R, as well as the robust competitors, that is, the SMDM-estimator, the CVFS-estimator and the Composite τ -estimator. A very important consideration in real situations is the selection of an “optimum” value of α that applies to the given data set. We will consider the values α^* and $\bar{\alpha}$, derived from theoretical computations, since they are the suggested optimal values.

We consider an orthodontic study conducted by Potthoff and Roy (1964) where the distance (in millimeters) between the pituitary and the pterygomaxillary fissure has been measured on 16 boys and 11 girls at 8, 10, 12, and 14 years. The data set is available as part of the R package `n.lme` (Pineiro et al. 2022). Figure 6 displays the measurements for each individual in the study together with the least-squares fit of the simple linear regression model.

From Fig. 6, it is possible to see that the data set possibly contains some outliers. In particular, the measurements for subject M09 have more variability around the fitted line with two possible within-subject outliers and the slope for subject M13 is larger than the others indicating a possible outlier at the level of random effects. Finally, subject M10 could be also considered as outlying observation for the large distance values since the first measurement. Overall, the intercept and the slope seem to vary with the subject, and the responses for the girls show less variation around the fitted lines than for boys.

According to the mentioned features, the orthodontic distance growth with respect to age can be modeled using the linear mixed model of the form

$$y_{ij} = \beta_0 + \beta_1 I_i(F) + (\beta_2 + \beta_3 I_i(F)) t_j + u_{i1} + u_{i2} t_j + \epsilon_{ij}$$

for $i = 1, \dots, 27$ and $j = 1, \dots, 4$, where y_{ij} denotes the distance for subject i at age t_j ; β_0 and β_1 represent the fixed effect intercept for boys and girls, respectively; β_2 and β_3 represent the fixed effect slope for boys and girls, respectively; $I_i(F)$ indicates an indicator function for the girls group; $(\mathcal{U}_1, \mathcal{U}_2)$ is the vector of random terms and (u_{i1}, u_{i2}) is the vector of realizations for subject i ; ϵ_{ij} is the error term. Notice that in this example, the random effect for the age is nested into the random effect for the subject level, so a possible correlation between them has to be considered. According to this, the variance components to be estimated are σ_j^2 , $j = 0, 1, 2$ as previously described and the covariance σ_{12} between the random variables \mathcal{U}_1 and \mathcal{U}_2 .

Table 4 and Table 5 report the estimates for the fixed terms and for the variance components, respectively, obtained by the MLE, the SMDM-estimator, the CVFS-estimator, the Composite τ -estimator and MDPDE for $\alpha = \alpha^* = 0.2$ and $\alpha = \bar{\alpha} = 0.5$. Furthermore, we computed the maximum likelihood estimates for the reduced data set obtained removing the observations M09, M10 and M13, indicated

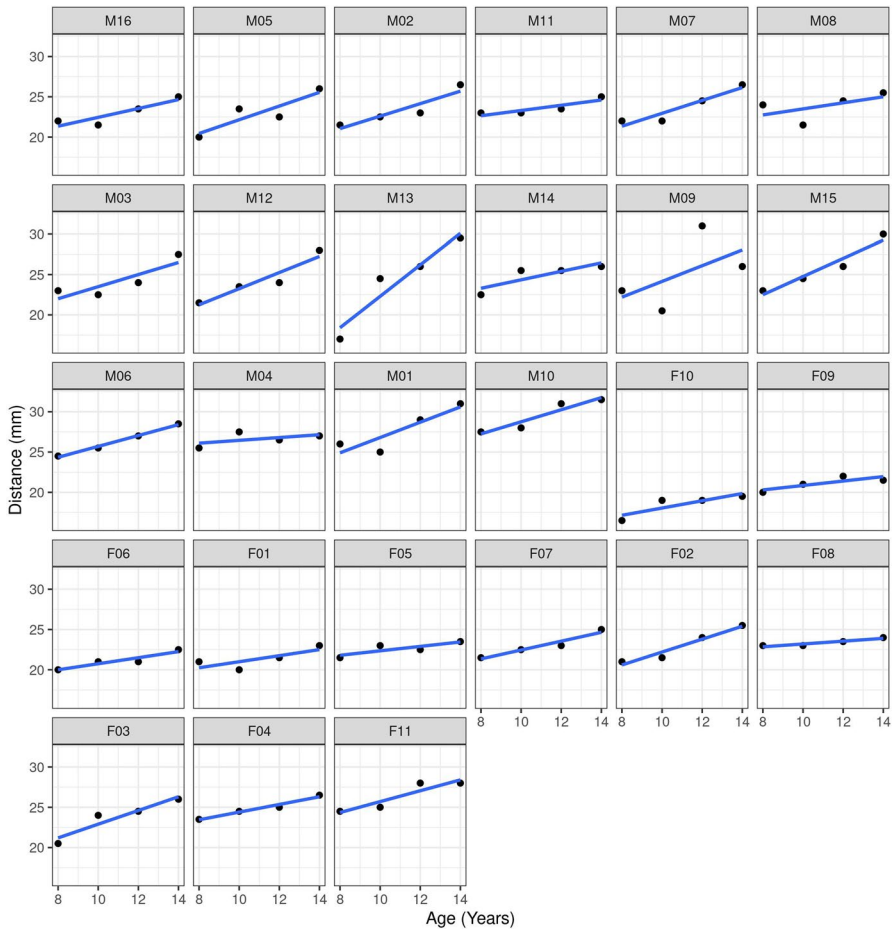


Fig. 6 Orthodontic growth patterns in 11 girls (F) and 16 boys (M) from 8 to 14 years of age. Blue lines correspond to the individual least-squares fit of the simple linear regression model

Table 4 Orthodont data set. Estimates of fixed effects parameters obtained by the MLE (with and without the outlying observations), the SMDM-estimator, the CVFS-estimator, the Composite τ -estimator and MDPDE for $\alpha = \alpha^* = 0.2$ and $\alpha = \bar{\alpha} = 0.5$

	MLE	SMDM	CSVF	Comp. τ	MDPD-0.2	MDPD-0.5	MLE (without)
$\hat{\beta}_0$	16.34	16.91	16.30	17.43	17.20	17.09	17.15
$\hat{\beta}_1$	1.03	0.53	1.32	0.08	0.35	0.74	0.22
$\hat{\beta}_2$	0.78	0.71	0.72	0.67	0.69	0.68	0.68
$\hat{\beta}_3$	-0.30	-0.23	-0.22	-0.20	-0.23	-0.24	-0.20

Table 5 Orthodont data set. Estimates of variance components parameters obtained by the MLE (with and without the outlying observations), the SMDM-estimator, the CVFS-estimator, the Composite τ -estimator and MDPDE for $\alpha = \alpha^* = 0.2$ and $\alpha = \bar{\alpha} = 0.5$

	MLE	SMDM	CVSF	Comp. τ	MDPD-0.2	MDPD-0.5	MLE (without)
$\hat{\sigma}_0^2$	1.72	1.21	1.41	1.09	0.93	0.95	0.88
$\hat{\sigma}_1^2$	5.79	1.08	0.46	2.00	3.51	3.06	3.35
$\hat{\sigma}_{12}$	-0.29	0.03	-0.00	0.02	-0.10	-0.09	-0.13
$\hat{\sigma}_2^2$	0.03	0.00	0.02	0.01	0.02	0.02	0.02

as outliers. The estimates of β_2 and β_3 parameters are quite similar among the considered estimators. The fixed effect parameters differ mainly for the intercept β_0 and β_1 , which indicates the effect of gender on the response y_i . Robust estimators show values closer to zero than the standard MLE, indicating that the orthodontic distance is less affected by gender. This can be explained by the fact that MLE is sensitive to the presence of outlying observations obtaining a bigger value. Indeed, when such observations are removed, the maximum likelihood estimates are similar to MDPDE. The main differences on the estimation of the random effects terms are both in size (error variance component) and shape (correlation components). The MDPDEs assign less variance to the error term, and the other estimates are in general smaller. As before, the estimates given by the MLE and MDPDE are very close when the outliers are removed. Notice that the SMDM- and CVFS-estimates are quite different with respect to the others, especially for the variance components parameters.

Figure 7 shows the QQ-plots of the estimates of the random terms u_{ij} for $i = 1, \dots, 27$ and $j = 1, 2$, obtained by the MLE, the SMDM-estimator, the CVFS-estimator, the Composite τ -estimator and MDPDE for $\alpha = \alpha^* = 0.2$ and $\alpha = \bar{\alpha} = 0.5$. The QQ-plots seem to show some structure and confirm the larger

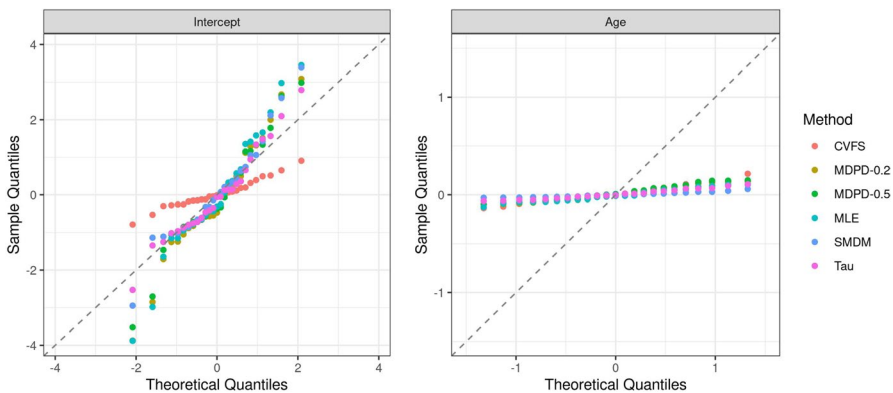


Fig. 7 Orthodont data set. QQ-plots of the random effects estimated by the MLE, the SMDM-estimator, the CVFS-estimator, the Composite τ -estimator and MDPDE for $\alpha = \alpha^* = 0.2$ and $\alpha = \bar{\alpha} = 0.5$

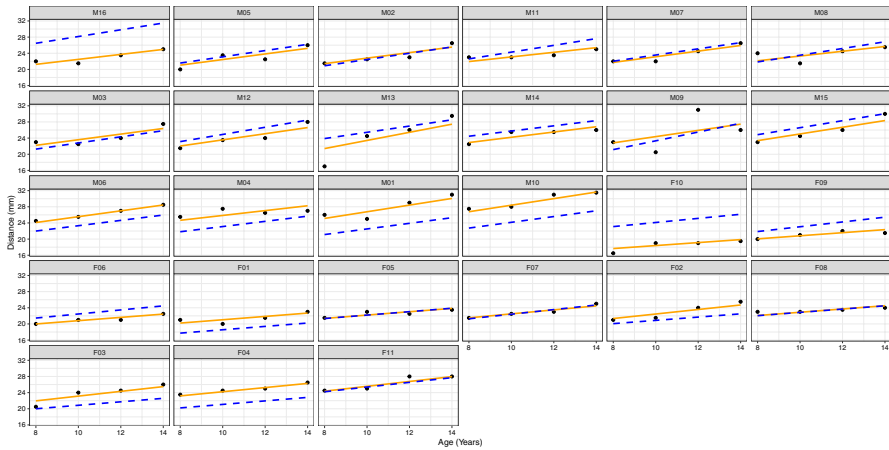


Fig. 8 Orthodontic growth patterns in 11 girls (F) and 16 boys (M) from 8 to 14 years of age. The fitted lines given by the MLE and MDPDE for $\alpha = 0.2$ are represented in blue and orange, respectively

variances estimated by MLE then those given by MDPDE; however, notice that the sample size ($n = 27$) is quite low. Figure 13a and 13b, in Section SM-7 of the Supplementary Materials, shows the QQ-plots of the random effects estimated by different methods separately.

Finally, Fig. 8 shows the fitted lines given by fixed effects estimates and random effects predictions for the MLE and MDPDE with $\alpha = 0.2$, represented in blue and orange, respectively. It is easy to see that the MDPDE achieves a better fit. In particular, notice the case of Subject M10 previously indicated as possible outlier.

Notice that for the simulation study in Sect. 5.2, the predictions of the response variable have been computed as $\hat{y} = X\hat{\beta}$, without considering the random effect estimates, since we were interested in evaluating the performance with respect to fixed effects estimation. In practical situations, such as the one presented in this section, when a linear mixed model is considered, both fixed effects and random effects estimates must be considered, i.e., $\hat{y}_i = X_i\hat{\beta} + Z_i\hat{u}_i$.

7 Conclusions

In this paper, we have developed an estimator based on the density power divergences to deal with the robustness issues in the linear mixed model setup. We demonstrated that the desirable asymptotic properties of the MDPDE, such as consistency and asymptotic normality, hold for the linear mixed model setup. In order to assess the robustness properties, the influence function and sensitivity measures of the estimator were computed. We found that the estimator is B-robust for $\alpha > 0$. From a practical point of view, the choice of the value of the tuning parameter α is fundamental in applications. The behavior of the sensitivity measures suggested two optimal values, denoted by α^* and $\bar{\alpha}$, depending on the dimension p , where the term “optimal” is in the sense of providing minimum sensitivity, and thus producing maximum robustness. The existence of such values is in contrast to the previous

knowledge about the parameter α . Indeed, it was shown that when α continues to increase beyond a certain value, we lose both robustness and efficiency.

The simulation study confirmed how the performance of the minimum density power divergence estimator changes with respect to α . Furthermore, the MDPDE outperforms the competitor estimators; indeed our approach leads to more resistant estimators in the presence of case-wise contamination. Finally, the application of our estimator to a real-life data set indicated that the MDPDE has similar results to the classical maximum likelihood estimator.

We feel that many important extensions of this work are necessary and can be potentially useful. So far, the MDPDE has been implemented only for balanced data (although the theory that we have developed is perfectly general). In future, we propose to extend the implementation to the more general case of groups with possibly different dimensions. The problem of testing of hypothesis also deserves a deeper look in the linear mixed models scenario.

Appendix A

A.1 Proof of Theorem 1

First let us note that under the setup of the linear mixed models introduced in Sect. 3, given a fixed $\alpha \geq 0$, for each i , the matrices $\mathbf{\Omega}_i$ and \mathbf{J}_i defined in Section SM-1 of Supplementary Material simplify to the forms

$$\mathbf{J}^{(i)} = \mathbb{E}_{g_i} [\nabla H_i(\mathbf{Y}_i, \boldsymbol{\theta})] = \begin{pmatrix} \mathbf{J}_{11}^{(i)} & 0 \\ 0 & \mathbf{J}_{22}^{(i)} \end{pmatrix},$$

$$\mathbf{\Omega}^{(i)} = \text{Var}_{g_i} (\nabla H_i(\mathbf{Y}_i, \boldsymbol{\theta})) = \begin{pmatrix} \mathbf{\Omega}_{11}^{(i)} & 0 \\ 0 & \mathbf{\Omega}_{22}^{(i)} \end{pmatrix},$$

where

$$\mathbf{J}_{11}^{(i)} = \frac{4\mathbf{X}_i^\top \mathbf{V}_i^{-1} \mathbf{X}_i}{(1 + \alpha)^{\frac{n_i}{2} + 1}},$$

(j, k) -th element of $\mathbf{J}_{22}^{(i)} = \frac{T(\mathbf{V}_i^{-1} \mathbf{U}_{ij}, \mathbf{V}_i^{-1} \mathbf{U}_{ik})}{(1 + \alpha)^{\frac{n_i}{2} + 2}},$

$$\mathbf{\Omega}_{11}^{(i)} = \frac{4\mathbf{X}_i^\top \mathbf{V}_i^{-1} \mathbf{X}_i}{(1 + 2\alpha)^{\frac{n_i}{2} + 1}},$$

(j, k) -th element of $\mathbf{\Omega}_{22}^{(i)} = \frac{T(4, \mathbf{V}_i^{-1} \mathbf{U}_{ij}, \mathbf{V}_i^{-1} \mathbf{U}_{ik})}{(1 + 2\alpha)^{\frac{n_i}{2} + 2}} - \frac{\alpha^2 \text{Tr}(\mathbf{V}_i^{-1} \mathbf{U}_{ik}) \text{Tr}(\mathbf{V}_i^{-1} \mathbf{U}_{ij})}{(1 + \alpha)^{n_i + 2}},$

for $j, k \in \{0, \dots, r\}$, with $T(c, \mathbf{A}, \mathbf{B}) = c\alpha^2 \text{Tr}(\mathbf{A}) \text{Tr}(\mathbf{B}) + 2\text{Tr}(\mathbf{A}\mathbf{B})$ for general matrices \mathbf{A}, \mathbf{B} and a constant c ($c = 1$ if not specified). Finally, put

$$\Psi_n = \frac{1}{n} \sum_{i=1}^n \frac{\eta_{i\alpha}}{4} \mathbf{J}^{(i)} \quad \text{and} \quad \Omega_n = \frac{1}{n} \sum_{i=1}^n \frac{\eta_{i\alpha}^2}{4} \Omega^{(i)}.$$

Also, we assumed that the true data generating density belongs to the model family. Then, the proof of our Theorem 1 is immediate from the results stated by Ghosh and Basu (2013), provided we can show that the general Assumptions (A1)–(A7) are satisfied under the our assumed Conditions (MM1)–(MM4) for the special case of LMMs.

Now, the assumption that the true data generating density belongs to the model family, together with that the model density is normal with mean $\mathbf{X}_i\boldsymbol{\beta}$ and variance matrix \mathbf{V}_i , ensures that Assumptions (A1)–(A3) are directly satisfied. Also, Assumption (A4) follows from Condition (MM1).

Next, we prove Equation (1) of Assumption (A6). For any $j = 1, \dots, p$, the j -th partial derivative with respect to $\boldsymbol{\beta}$ is given by

$$\nabla_j H_i(\mathbf{Y}_i, \boldsymbol{\theta}) = -\alpha \left(1 + \frac{1}{\alpha} \right) \eta_{i\alpha} e^{-\frac{\alpha}{2}(\mathbf{Y}_i - \mathbf{X}_i\boldsymbol{\beta})^\top \mathbf{V}_i^{-1}(\mathbf{Y}_i - \mathbf{X}_i\boldsymbol{\beta})} \mathbf{X}_{ij}^\top \mathbf{V}_i^{-1}(\mathbf{Y}_i - \mathbf{X}_i\boldsymbol{\beta}).$$

Then, considering $\mathbf{Z}_i = \mathbf{V}_i^{-\frac{1}{2}}(\mathbf{Y}_i - \mathbf{X}_i\boldsymbol{\beta})$,

$$\begin{aligned} & \frac{(1 + \alpha)}{n} \sum_{i=1}^n E_{g_i} \left[\left| \eta_{i\alpha} e^{-\frac{\alpha}{2}(\mathbf{Y}_i - \mathbf{X}_i\boldsymbol{\beta})^\top \mathbf{V}_i^{-1}(\mathbf{Y}_i - \mathbf{X}_i\boldsymbol{\beta})} \mathbf{X}_{ij}^\top \mathbf{V}_i^{-1}(\mathbf{Y}_i - \mathbf{X}_i\boldsymbol{\beta}) \right| \right. \\ & \quad \left. \times uadl \left(\left| \eta_{i\alpha} e^{-\frac{\alpha}{2}(\mathbf{Y}_i - \mathbf{X}_i\boldsymbol{\beta})^\top \mathbf{V}_i^{-1}(\mathbf{Y}_i - \mathbf{X}_i\boldsymbol{\beta})} \mathbf{X}_{ij}^\top \mathbf{V}_i^{-1}(\mathbf{Y}_i - \mathbf{X}_i\boldsymbol{\beta}) \right| > \frac{N}{1 + \alpha} \right) \right] \\ & \leq \frac{(1 + \alpha)}{n} \sum_{i=1}^n \eta_{i\alpha} |\mathbf{X}_{ij}^\top \mathbf{V}_i^{-\frac{1}{2}}| E_{g_i} \left[e^{-\frac{\alpha}{2} \mathbf{Z}_i^\top \mathbf{Z}_i} |\mathbf{Z}_i| \right. \\ & \quad \left. \times I \left(e^{-\frac{\alpha}{2} \mathbf{Z}_i^\top \mathbf{Z}_i} |\mathbf{Z}_i| > \frac{N}{(1 + \alpha)(\sup_{n>1} \eta_{i\alpha})(\sup_{n>1} \max_{1 \leq i \leq n} |\mathbf{X}_{ij}^\top \mathbf{V}_i^{-\frac{1}{2}}|)} \right) \right] \\ & = E_1 \left[e^{-\frac{\alpha}{2} \mathbf{Z}_1^\top \mathbf{Z}_1} |\mathbf{Z}_1| I \left(e^{-\frac{\alpha}{2} \mathbf{Z}_1^\top \mathbf{Z}_1} |\mathbf{Z}_1| > \frac{N}{(1 + \alpha)(\sup_{n>1} \eta_{i\alpha})(\sup_{n>1} \max_{1 \leq i \leq n} |\mathbf{X}_{ij}^\top \mathbf{V}_i^{-\frac{1}{2}}|)} \right) \right] \\ & \quad \times \left(\frac{1 + \alpha}{n} \sum_{i=1}^n \eta_{i\alpha} |\mathbf{X}_{ij}^\top \mathbf{V}_i^{-\frac{1}{2}}| \right) \end{aligned}$$

Since the $\sup_{n>1} \max_{1 \leq i \leq n} |\mathbf{X}_{ij}^\top \mathbf{V}_i^{-\frac{1}{2}}| = \mathcal{O}(1)$ by Assumption (MM2) and $\sup_{n>1} \eta_{i\alpha}$ is bounded thanks to the boundness of $|\mathbf{V}_i|$ in (MM3), by the dominated convergence theorem (DCT), we have

$$\lim_{N \rightarrow \infty} E_1 \left[e^{-\frac{\alpha}{2} \mathbf{Z}_1^\top \mathbf{Z}_1} |\mathbf{Z}_1| \left| I \left(e^{-\frac{\alpha}{2} \mathbf{Z}_1^\top \mathbf{Z}_1} |\mathbf{Z}_1| > \frac{N}{(1 + \alpha)(\sup_{n>1} \eta_{i\alpha})(\sup_{n>1} \max_{1 \leq i \leq n} |\mathbf{X}_{ij}^\top \mathbf{V}_i^{-\frac{1}{2}}|)} \right) \right] = 0.$$

Also

$$\sup_{n>1} \left(\frac{1}{n} \sum_{i=1}^n \eta_{i\alpha} |\mathbf{X}_{ij}^\top \mathbf{V}_i^{-\frac{1}{2}}| \right) \leq (\sup_{n>1} \eta_{i\alpha})(\sup_{n>1} \max_{1 \leq i \leq n} |\mathbf{X}_{ij}^\top \mathbf{V}_i^{-\frac{1}{2}}|) = \mathcal{O}(1)$$

and this follows for all $j = 1, \dots, p$. On the other hand, consider the partial derivative with respect to σ_j^2 , $j = 0, 1, \dots, r$, given in Equation (10). Hence, denoting with

$\mathbf{Z}_i = \mathbf{V}_i^{-\frac{1}{2}}(\mathbf{Y}_i - \mathbf{X}_i \boldsymbol{\beta})$, we have

$$\begin{aligned} & \frac{1}{n} \sum_{i=1}^n E_{g_i} \left[|\nabla_{\sigma_j^2} H_i(\mathbf{Y}_i, \boldsymbol{\theta})| \times I(|\nabla_{\sigma_j^2} H_i(\mathbf{Y}_i, \boldsymbol{\theta})| > N) \right] \\ & \leq \frac{\alpha}{2n} \sum_{i=1}^n \frac{\eta_{i\alpha} \text{Tr}(\mathbf{V}_i^{-1} \mathbf{U}_{ij})}{(1 + \alpha)^{\frac{\eta_i}{2}}} E_{g_i} [I(|\nabla_{\sigma_j^2} H_i(\mathbf{Y}_i, \boldsymbol{\theta})| > N)] \\ & \quad + \frac{1 + \alpha}{2n} \sum_{i=1}^n \eta_{i\alpha} \text{Tr}(\mathbf{V}_i^{-1} \mathbf{U}_{ij}) E_{g_i} [e^{-\frac{\alpha}{2} \mathbf{Z}_i^\top \mathbf{Z}_i} I(|\nabla_{\sigma_j^2} H_i(\mathbf{Y}_i, \boldsymbol{\theta})| > N)] \\ & \quad + \frac{1 + \alpha}{2n} \sum_{i=1}^n \eta_{i\alpha} E_{g_i} [|\mathbf{Z}_i^\top \mathbf{V}_i^{-1} \mathbf{U}_{ij} \mathbf{V}_i^{-1} \mathbf{Z}_i| e^{-\frac{\alpha}{2} \mathbf{Z}_i^\top \mathbf{V}_i^{-1} \mathbf{Z}_i} I(|\nabla_{\sigma_j^2} H_i(\mathbf{Y}_i, \boldsymbol{\theta})| > N)] \\ & = \left(\frac{\alpha}{2n} \sum_{i=1}^n \frac{\eta_{i\alpha} \text{Tr}(\mathbf{V}_i^{-1} \mathbf{U}_{ij})}{(1 + \alpha)^{\frac{\eta_i}{2}}} \right) E_1 [I(|\nabla_{\sigma_j^2} H_i(\mathbf{Y}_i, \boldsymbol{\theta})| > N)] \\ & \quad + \left(\frac{1 + \alpha}{2n} \sum_{i=1}^n \eta_{i\alpha} \text{Tr}(\mathbf{V}_i^{-1} \mathbf{U}_{ij}) \right) E_1 [e^{-\frac{\alpha}{2} \mathbf{Z}_i^\top \mathbf{Z}_i} I(|\nabla_{\sigma_j^2} H_i(\mathbf{Y}_i, \boldsymbol{\theta})| > N)] + \\ & \quad + \frac{1 + \alpha}{2n} \sum_{i=1}^n \eta_{i\alpha} E_{g_i} [|\mathbf{Z}_i^\top \mathbf{V}_i^{-1} \mathbf{U}_{ij} \mathbf{V}_i^{-1} \mathbf{Z}_i| e^{-\frac{\alpha}{2} \mathbf{Z}_i^\top \mathbf{V}_i^{-1} \mathbf{Z}_i} I(|\nabla_{\sigma_j^2} H_i(\mathbf{Y}_i, \boldsymbol{\theta})| > N)] \end{aligned}$$

where in the last term $\mathbf{Z}_i = (\mathbf{Y}_i - \mathbf{X}_i \boldsymbol{\beta})$. Note that

$$E_{g_i} [|\mathbf{Z}_i^\top \mathbf{V}_i^{-1} \mathbf{U}_{ij} \mathbf{V}_i^{-1} \mathbf{Z}_i| e^{-\frac{\alpha}{2} \mathbf{Z}_i^\top \mathbf{V}_i^{-1} \mathbf{Z}_i}] = \frac{\text{Tr}(\mathbf{V}_i^{-1} \mathbf{U}_{ij})}{(1 + \alpha)^{1 + \frac{\eta_i}{2}}}$$

; hence, by Equation (15),

$$E_{g_i} \left[|\mathbf{Z}_i^\top \mathbf{V}_i^{-1} \mathbf{U}_{ij} \mathbf{V}_i^{-1} \mathbf{Z}_i| e^{-\frac{\alpha}{2} \mathbf{Z}_i^\top \mathbf{V}_i^{-1} \mathbf{Z}_i} I(|\nabla_{\sigma_j^2} H_i(\mathbf{Y}_i, \boldsymbol{\theta})| > N) \right] \rightarrow 0, \quad \text{as } N \rightarrow \infty.$$

The expectation in the first two terms goes to zero as $N \rightarrow \infty$ by the DCT, as before, and the sums are bounded by Equation (15) in condition (MM3). This holds for all $j = 0, \dots, r$.

Finally, Assumptions (A5), (A7) and Equation (2) similarly hold using Equations (14) and (16), (17), (13) and (15), respectively.

Supplementary Information The online version contains supplementary material available at <https://doi.org/10.1007/s10182-023-00473-z>.

Acknowledgements The research of AG is supported by an INSPIRE Faculty Research Grant from Department of Science and Technology, Government of India, and an internal research grant from Indian Statistical Institute, India.

Author Contributions G.S. designed the methodology, developed the implementation of the proposed method and prepared the visualization of the results as well as the preparation of the original presented work and of the revised version. A.G. formulated the initial idea and the research goal and aims and validated the theoretical and experimental research output. A.B. is responsible of coordinating the research activity. C.A. supervised the planning and the execution of the project and validated the implementation of the method.

Funding Open access funding provided by Università degli Studi di Trento within the CRUI-CARE Agreement.

Declarations

Conflict of interest We have no conflicts of interest to disclose.

Open Access This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

References

- Agostinelli, C., Yohai, V.J.: *robustvarComp: Robust Estimation for Variance Component Models*. (2019). R package version 0.1-6
- Agostinelli, C., Yohai, V.J.: Composite robust estimators for linear mixed models. *J. Am. Stat. Assoc.* **111**(516), 1764–1774 (2016)
- Basu, A., Harris, I.R., Hjort, N., Jones, M.C.: Robust and efficient estimation by minimizing a density power divergence. *Biometrika* **85**(3), 549–559 (1998)
- Basu, A., Park, C., Shioya, H.: *Statistical Inference: The Minimum Distance Approach*. CRC Press, Chapman and Hall (2011)
- Castilla, E., Ghosh, A., Martin, N., Pardo, L.: New robust statistical procedures for the polytomous logistic regression models. *Biometrics* **74**(4), 1282–1291 (2018)

- Castilla, E., Ghosh, A., Martin, N., Pardo, L.: Robust semiparametric inference for polytomous logistic regression with complex survey design. *Adv. Data Anal. Classificat.* **15**(3), 701–734 (2021). <https://doi.org/10.1007/s11634-020-00430-7>
- Christensen, R.: Mixed models and variance components. In: *Plane Answers to Complex Questions: The Theory of Linear Models*, pp. 291–331. Springer, New York, NY (2011)
- Copt, S., Victoria-Feser, M.P.: High breakdown inference in the mixed linear model. *J. Am. Stat. Assoc.* **101**, 292–300 (2006)
- Ghosh, A.: Robust inference under the beta regression model with application to health care studies. *Stat. Methods Med. Res.* **28**(3), 871–888 (2019)
- Ghosh, A., Basu, A.: Robust estimation for independent non-homogeneous observations using density power divergence with applications to linear regression. *Electr. J. Stat.* **7**, 2420–2456 (2013)
- Ghosh, A., Basu, A.: Robust estimation in generalized linear models: The density power divergence approach. *TEST* **25**, 269–290 (2016)
- Ghosh, A., Basu, A.: Robust and efficient estimation in the parametric proportional hazards model under random censoring. *Stat. Med.* **38**(27), 5283–5299 (2019)
- Hampel, F.R.: *Contributions to the Theory of Robust Estimation*. University of California, Berkeley (1968)
- Hampel, F.R.: The influence curve and its role in robust estimation. *J. Am. Stat. Assoc.* **69**(346), 383–393 (1974)
- Huggins, R.M.: On the robust analysis of variance components models for pedigree data. *Aust. J. Stat.* **35**(1), 43–57 (1993)
- Huggins, R.M.: A robust approach to the analysis of repeated measures. *Biometrics* **49**(3), 715–720 (1993)
- Huggins, R.M., Staudte, R.G.: Variance components models for dependent cell populations. *J. Am. Stat. Assoc.* **89**(425), 19–29 (1994)
- Koller, M.: Robust estimation of linear mixed models. PhD thesis, ETH Zürich (2013)
- Koller, K.: *robustlmm: An R package for robust estimation of linear mixed-effects models*. *J. Stat. Softw.* **75**(6), 1–24 (2016)
- Lange, K.L., Little, R.J.A., Taylor, J.M.G.: Robust statistical modeling using the t distribution. *J. Am. Stat. Assoc.* **84**(408), 881–896 (1989)
- McCulloch, C.E., Searle, S.R.: *Generalized, Linear, and Mixed Models*. John Wiley & Sons, Wiley Series in Probability and Statistics (2001)
- Pinheiro, J., Bates, D., R Core Team: *Nlme: Linear and Nonlinear Mixed Effects Models*. (2022). R package version 3.1-157. <https://CRAN.R-project.org/package=nlme>
- Pinheiro, J.C., Liu, C., Wu, Y.N.: Efficient algorithms for robust estimation in linear mixed-effects models using the multivariate t distribution. *J. Comput. Graph. Stat.* **10**(2), 249–276 (2001)
- Potthoff, R.F., Roy, S.N.: A generalized multivariate analysis of variance model useful especially for growth curve problems. *Biometrika* **51**(3/4), 313–326 (1964)
- R Core Team: *R: A Language and environment for statistical computing*. R Foundation for Statistical Computing, Vienna, Austria (2019). R Foundation for Statistical Computing
- Richardson, A.M.: Bounded influence estimation in the mixed linear model. *J. Am. Stat. Assoc.* **92**(437), 154–161 (1997)
- Richardson, A.M., Welsh, A.H.: Robust restricted maximum likelihood in mixed linear models. *Biometrics* **51**(4), 1429–1439 (1995)
- Sinha, S.K.: Robust analysis of generalized linear mixed models. *J. Am. Stat. Assoc.* **99**(466), 451–460 (2004)
- Stahel, W.A., Welsh, A.: Approaches to robust estimation in the simplest variance components model. *J. Stat. Plan. Infer.* **57**(2), 295–319 (1994)
- Welsh, A.H., Richardson, A.M.: 13 approaches to the robust estimation of mixed models. *Handbook Stat.* **15**, 343–384 (1997)
- Yau, K.K.W., Kuk, A.Y.C.: Robust estimation in generalized linear mixed models. *J. Royal Stat. Soc.* **64**(1), 101–117 (2002)

Authors and Affiliations

Giovanni Saraceno^{1,3}  · **Abhik Ghosh**² · **Ayanendranath Basu**² · **Claudio Agostinelli**¹

Abhik Ghosh
abhik.ghosh@isical.ac.in

Ayanendranath Basu
ayanbasu@isical.ac.in

Claudio Agostinelli
claudio.agostinelli@unitn.it

- ¹ Department of Mathematics, University of Trento, Via Sommarive 14, 38123 Trento, Italy
- ² Interdisciplinary Statistical Research Unit, Indian Statistical Institute, 203 Barrackpore Trunk Road Kolkata, Kolkata 700 108, India
- ³ Department of Biostatistics, University at Buffalo, 811 Kimball Tower, Buffalo, NY 14214, USA