Made available by Hasselt University Library in https://documentserver.uhasselt.be

Estimating the reliability of repeatedly measured endpoints based on linear mixed-effects models. A tutorial Peer-reviewed author version

VAN DER ELST, Wim; MOLENBERGHS, Geert; Hilgers, Ralf-Dieter; VERBEKE, Geert & Heussen, Nicole (2016) Estimating the reliability of repeatedly measured endpoints based on linear mixed-effects models. A tutorial. In: PHARMACEUTICAL STATISTICS, 15(6), p. 486-493.

DOI: 10.1002/pst.1787 Handle: http://hdl.handle.net/1942/23064

Estimating the reliability of repeatedly measured endpoints based on linear mixed-effects models. A tutorial

Wim Van der Elst,¹ Geert Molenberghs^{1,2}, Ralf-Dieter Hilgers,³ Geert Verbeke^{1,2} & Nicole Heussen³

Abstract

There are various settings in which researchers are interested in the assessment of the correlation between repeated measurements that are taken *within* the same subject (i.e., reliability). For example, the same rating scale may be used to assess the symptom severity of the same patients by multiple physicians, or the same outcome may be measured repeatedly over time in the same patients.

Reliability can be estimated in various ways, e.g., using the classical Pearson correlation or the intra-class correlation in clustered data. However, contemporary data often have a complex structure that goes well beyond the restrictive assumptions that are needed with the more conventional methods to estimate reliability.

In the current paper, we propose a general and flexible modeling approach that allows for the derivation of reliability estimates, standard errors, and confidence intervals – appropriately taking hierarchies and covariates in the data into account. Our methodology is developed for continuous outcomes together with covariates of an arbitrary type.

The methodology is illustrated in a case study, and a Web Appendix is provided which details the computations using the R package *CorrMixed* and the SAS software.

Keywords: within-cluster correlation; test-retest reliability; intra-class correlation

¹I-BioStat, Universiteit Hasselt, B-3590 Diepenbeek, Belgium.

²I-BioStat, Katholieke Universiteit Leuven, B-3000 Leuven, Belgium.

³Department of Medical Statistics, RWTH Aachen University, Aachen, Germany

1 Introduction

Reliability essentially refers to the reproducibility (or, predictability) of outcomes that are repeatedly measured *within* the same individuals. In particular, this metric quantifies the extent to which a repetition of a measurement
under the same general conditions leads to the same result.

Conventional methods to estimate reliability The concept of relia-5 bility is grounded in the so-called classical test theory [1]. In this paradigm, 6 the outcome of a measurement procedure is modeled as $X = \tau + \varepsilon$, where 7 X is the observed score of a subject, τ is the unobserved (latent) true score 8 of this person, and ε is the measurement error. In classical test theory, it 9 is assumed (i) that the measurement errors are mutually uncorrelated, and 10 (ii) that the measurement errors are uncorrelated with the true scores. Un-11 der these assumptions, $\operatorname{Var}(X) = \operatorname{Var}(\tau) + \operatorname{Var}(\varepsilon)$ and the reliability of the 12 measurement (R) is defined as 13

$$R = \frac{\operatorname{Var}(\tau)}{\operatorname{Var}(X)} = \frac{\operatorname{Var}(\tau)}{\operatorname{Var}(\tau) + \operatorname{Var}(\varepsilon)}.$$
(1)

Eq. (1) is intuitively appealing because it defines reliability as the fraction of 14 the observed test score variance that is attributable to the true score variance. 15 If a test is perfectly reliable, the true score and observed score variances are 16 equal and thus R = 1. Unfortunately, reliability cannot be directly estimated 17 based on Eq. (1) because τ cannot be observed. Instead, reliability will have 18 to be estimated indirectly. A classical solution to the problem is to introduce 19 the concept of *parallel tests* [2]. Parallel tests are tests that have the same 20 true score for each subject and equal error variances. For example, suppose 21 that we have two measurements X_1 and X_2 for the same subjects that are 22 assessed at two instances of time with a short lag (such that τ does not 23 change), or that are obtained from two raters at the same point in time. 24 Then $X_1 = \tau + \varepsilon_1$ and $X_2 = \tau + \varepsilon_2$ with $\operatorname{Var}(X_1) = \operatorname{Var}(X_2) = \operatorname{Var}(X)$ and 25 $\operatorname{Var}(\varepsilon_1) = \operatorname{Var}(\varepsilon_2) = \operatorname{Var}(\varepsilon)$, i.e., X_1 and X_2 are parallel measurements. The 26 covariance of the two measurements then equals 27

$$Cov(X_1, X_2) = Cov(\tau + \varepsilon_1, \tau + \varepsilon_2)$$

= Var(\tau) + Cov(\tau, \varepsilon_1) + Cov(\tau, \varepsilon_2) + Cov(\varepsilon_1, \varepsilon_2)
= Var(\tau),

and the correlation between X_1 and X_2 can be written as

$$\operatorname{Corr}(X_1, X_2) = \frac{\operatorname{Cov}(X_1, X_2)}{\sqrt{\operatorname{Var}(X_1)}\sqrt{\operatorname{Var}(X_2)}} = \frac{\operatorname{Var}(\tau)}{\operatorname{Var}(\tau) + \operatorname{Var}(\varepsilon)} = R.$$
(2)

Limitations of the conventional methods Eq. (2) provides a conve-29 nient and straightforward way to compute reliability, but it is important 30 to stress that the assumption that the measurements are parallel is crucial. 31 This assumption is often violated in practice [3]. For example, it seems im-32 plausible to assume that patients in a clinical trial or in medical practice 33 do not exhibit a systematic change over time as a result of their treatment. 34 Another limitation of Eq. (2) is that only two measurements can be consid-35 ered, and these measurements should have the same test-retest interval for 36 all subjects. In practice, data may be available for more than two measure-37 ment moments and/or with different test-retest intervals. Further, the use of 38 Eq. (2) is less-than-ideal when data are missing, because subjects who have a 39 missing observation for either X_1 or X_2 are discarded from the analysis. This 40 approach does not only lead to a loss of information, but it also ignores the 41 missing data generating mechanism. Basically, to obtain unbiased estimates 42 for R using Eq. (2), the assumption that the data are missing completely at 43 random (MCAR) should be valid. This means that the missingness should 44 not depend on the observed or the unobserved outcomes [4, 5]. This is a 45 strong and often unrealistic assumption, e.g., in a clinical trial setting it is 46 conceivable that subjects who have lower scores at the first measurement in 47 time (poorer health) are more likely to drop out of the study at the second 48 measurement in time (missing value for X_2). 49

Importance of reliability It is important to carefully consider the relia-50 bility of a measurement procedure, for example in the context of designing 51 a clinical trial. Obviously, in particular in explorative or experimental small 52 population group studies, serial measurements are gathered to understand 53 the nature of the disease. However, unreliable measurement methods might 54 lead to serious misinterpretation of the disease process. Indeed, even the 55 most elegant study design will not overcome the damage that is caused by 56 the use of unreliable measurement procedures [6]. For example, biased sam-57 ple selection may occur when patients are selected based on an unreliable 58 measurement procedure, and the sample size that is required to detect an 59

28

important treatment difference (δ) may increase substantially when the out-60 come of interest is quantified using an unreliable measurement procedure. As 61 an illustration of the latter issue, consider a situation where a t-test is used 62 to evaluate the treatment effect on the primary endpoint in a clinical trial 63 with two treatment groups. When the measurement procedure that is used 64 to quantify the primary endpoint has perfect reliability (i.e., R = 1), the 65 required sample size to detect δ equals n^* . However, when this measurement 66 procedure has a less-than-perfect reliability (i.e., R < 1), the required sample 67 size becomes $n = \frac{n^*}{R}$ (for details, see [6]). Thus, for example, when R = 0.50, 68 the required sample size to detect δ doubles compared to what would have 69 been needed when R = 1. Clearly, an increase in the required sample size is 70 an issue in nearly all clinical studies (e.g., increased study duration and cost) 71 - and it may even make the conduct of the study infeasible (e.g., clinical 72 trials in rare diseases). 73

Aim and organization of the paper The main aim of the present pa-74 per is to illustrate how reliability can be estimated in a flexible way using 75 linear mixed-effects models (LMMs). As will be detailed below, LMMs can 76 separate the mean and the variance structures in the data – which allows 77 for relaxing the strong assumptions that are needed to apply the conven-78 tional methods to estimate reliability. Further, LMMs can deal with data 79 structures where different subjects have a different number of repeated mea-80 surements (2 or more) – which may or may not be regularly spaced. Finally, 81 LMMs are so-called likelihood based methods that provide valid results when 82 the missingness mechanism is missing at random (MAR) [7]. MAR means 83 that the missingness may depend on the observed outcomes (e.g., the first 84 measurement X_1) but not on unobserved outcomes. MAR is a substantially 85 less restrictive assumption than MCAR, and is thus more likely to hold in 86 practice 4. 87

The remainder of the paper is organized in the following way. In Section 88 2, a case study is introduced that will be used throughout this paper to 89 illustrate the methodology. In Section 3, an exploratory analysis of the case 90 study is conducted. In Section 4, the LMM-based approach to estimate 91 reliability is detailed. Section 5 discusses the results. A Web Appendix is 92 also provided in which additional materials are presented. In particular, it 93 details all the required computations using the newly developed R software 94 package CorrMixed and SAS. 95

$\mathbf{2}$ Case study 96

114

Pikkemaat et al. [8] performed an experiment where the cardiac output 97 and stroke volume of N = 14 pigs was changed by increasing positive end-98 expiratory pressure (PEEP) levels $(0, 5, 10, 15, 20, \text{ and } 25 \text{ cm } H_2 \text{O})$. The 90 number of times that a particular PEEP level was used varied from animal 100 to animal. For each PEEP level, stroke volume was measured by the contin-101 uous approximately normally distributed variable Electrical Impedance To-102 mography (ZSV). In each animal, four identical experiments were conducted 103 (referred to as Cycles 1 to 4). The number of repeated ZSV measurements 104 across PEEP levels and cycles in an animal ranged between 9 and 47. In the 105 analyses below, it is assumed that all the measurements are equally spaced. 106 Pikkemaat et al. [8] were interested in estimating the levels of association 107 between the repeatedly measured ZSV and SVTTD (transpulmonary ther-108 modilution) outcomes within an animal. As detailed in the Introduction, it 109 is also worthwhile to evaluate the reliability of these repeated measurements. 110 Such analyses (not considered in [8]) will be the focus of the current paper. 111 Given the complex design of the study, it is recommended to use a flexible 112 LMM-based technique to estimate reliability (see Section 4) - rather than 113 the conventional techniques that were discussed in the Introduction.

As noted above, the study included a total of 14 pigs. However, the data 115 of n = 2 animals could not be evaluated due to technical reasons and these 116 animals were thus excluded from the analyses. Further, there were n = 2117 animals who appeared to have a 'clinically deviating' profile (as judged by 118 the experimenters). These animals were kept in the current analyses, but a 119 sensitivity analysis showed that the estimated reliabilities were not substan-120 tially affected by the in- or exclusion of these animals (see Web Appendix 121 Part II). Note that the data for PEEP level 25 were included in the current 122 analysis, as well as in the Pikkemaat et al. [8] study, although they were not 123 explicitly mentioned in the latter. 124

3 Exploratory data analysis 125

Figure 1 shows the individual profiles (grev lines) of ZSV as a function of 126 measurement moment. As can be seen, there is substantial between- as well 127 as within-animal variability. Further, drop-out is substantial, i.e., there are 128 less observations at later measurement moments compared to earlier mea-129

surement moments. This is more clearly depicted in Figure 2, where the
number of available observations at each of the different measurement moments are shown.

Figure 1 also shows that the average evolution over time (solid black line) exhibits a rather complex shape that cannot be modeled in a straightforward way by using linear or quadratic polynomials. Therefore, it is useful to consider a more general family of parametric models that are based on so-called fractional polynomial functions [9].

- 138
- $_{139}$ » Figures 1 and 2 about here «

Fractional polynomials The idea is to fit regression models with m terms of the form t^p , where the exponents p are selected from a small predefined set S of both integer and non-integer values. The linear predictor for a fractional polynomial of order M for covariate t (here: measurement point in time) on the mean ZSV is then defined as:

$$\beta_0 + \sum_{m=1}^M \beta_m t^{p_m}.$$
(3)

Each power p_m is chosen from a restricted set, typically $S = \{-2, -1, -0.5, 0, 0.5, 1, 2, 3\}$. 145 Note that when M = 2 and $p_1 = p_2$, the linear predictor (3) becomes 146 $\beta_0 + \beta_1 t^{p_1} + \beta_2 t^{p_1} \log(t)$. Also, when p = 0, this is taken to refer to $\log(t)$ [9]. 147 In practice, all possible models of degree 1 to M are fitted. Thus for M = 1, 148 each of the 8 values of S are used for the predictor t^{p_1} , for M = 2 each of the 149 36 combinations of powers are used for the predictors t^{p_1} and t^{p_2} , and so on. 150 Subsequently, the 'best' fitting model is selected. This choice can be made 151 in an informal way (i) based on Akaike's Information Criterion (AIC, where 152 a lower value is indicative of a better model fit) and/or (ii) by graphically 153 evaluating the fit of the model with the observed data. The AIC adds the 154 number of model parameters as a penalty to the log likelihood of the model, 155 which may help to avoid over-fitting (even though one still may want to be 156 careful not to select an overly complex model, in particular when a large 157 number of candidate powers is considered). The main advantage of using 158 fractional polynomials (rather than regular polynomials) is that they allow 159 for a much more flexible parametrization, i.e., a large number of different 160 shapes of curves can be captured by even a relatively small M. 161

Application to the case study In the analysis of the case study, frac-162 tional polynomials of order M = 1 to M = 5 were considered using the 163 standard set $S = \{-2, -1, -0.5, 0, 0.5, 1, 2, 3\}$ for the powers p_m . Note 164 that it is possible to use a more extensive set of values for S if the original set 165 does not provide an adequate result, but the number of models that have to 166 be fitted (and thus also the required computational time) increases sharply 167 when the number of elements in S increases. For example, when the set S168 includes 8 elements (the standard set), a total of 792 fractional polynomials 169 of degree 5 can be made. However, when the set $S = \{-3, -2.75, ..., 3\}$ 170 is used (25 elements), a total of 118,755 fractional polynomials of degree 5 171 can be made. Similarly, M can be increased but this will again yield a sharp 172 increase in the number of models to be evaluated. 173

Thus, regression models that included linear predictors for fractional polynomials of order M = 1 to M = 5 (see Eq. (3)) were fitted to the data of the case study. Table 1 shows the powers p_m of the models of order 1 to 5 that had the lowest AIC values. As can be seen, the model with M = 3 had the lowest overall AIC value. Figure 3 shows the predicted mean ZSV as a function of measurement moment for this model.

Based on these results, the fractional polynomial of degree 3 was retained as the 'best' model for the subsequent analyses. Thus, in the LMM analyses detailed below, the relation between time of measurement t and the mean ZSV will be modeled as $\beta_1 t^2 + \beta_2 t^2 \log(t) + \beta_3 t^3$.

184

185 » Table 1 about here «

¹⁸⁶ » Figure 3 about here «

¹⁸⁷ 4 Estimating reliability using mixed-effects mod ¹⁸⁸ els

In this section the reliability of the ZSV will be estimated using a flexible approach that is based on LMMs. The LMM is briefly introduced in Section 4.1 (for more details, see e.g., [7, 10, 11]), and the LMM-based approach to estimate reliability is applied to the case study in Section 4.2. For conciseness, in the latter section only a summary of the main results is given and no reference to software tools that can be used to obtain the results is made. However, full details can be found in the Web Appendix Parts I–V.

¹⁹⁶ 4.1 The linear mixed-effects model

¹⁹⁷ A LMM can be written as

$$\mathbf{Y}_i = \mathbf{X}_i \boldsymbol{\beta} + \mathbf{Z}_i \mathbf{b}_i + \boldsymbol{\varepsilon}_i, \tag{4}$$

where \mathbf{Y}_i is the response vector for subject *i* (with i = 1, 2, ..., n subjects 198 in the study), \mathbf{X}_i and \mathbf{Z}_i are the known design matrices for the fixed and 199 random effects, β is the vector that contains the fixed effects, \mathbf{b}_i is the vector 200 that contains the random effects, and ε_i is the vector that contains the mea-201 surement error (with $\mathbf{b}_i \sim N(\mathbf{0}, \mathbf{D})$ and $\boldsymbol{\varepsilon}_i \sim N(\mathbf{0}, \boldsymbol{\Sigma}_i)$, where \mathbf{D} and $\boldsymbol{\Sigma}_i$ 202 are general variance-covariance matrices). Model (4) thus assumes that the 203 vector of repeated measurements for each subject follows a linear regression 204 model where some of the parameters are population-specific (that is, param-205 eters that are the same for all subjects in the population; the fixed effects) 206 and other parameters are subject-specific (that is, parameters that differ for 207 all subjects; the random effects). 208

The residual component ε_i is often further decomposed as $\varepsilon_i = \varepsilon_{(1)i} + \varepsilon_i$ 209 $\boldsymbol{\varepsilon}_{(2)i}$. Here, $\boldsymbol{\varepsilon}_{(2)i}$ is a component of serial correlation and $\boldsymbol{\varepsilon}_{(1)i}$ is a component 210 of measurement error. Serial correlation results from the fact that within 211 a subject, the (residuals of) observations that are closer in time are often 212 'more similar' (i.e., more strongly correlated) than observations that are more 213 distant in time. It is assumed that $\boldsymbol{\varepsilon}_{(1)i} \sim N(\mathbf{0}, \sigma^2 \mathbf{I}_{n_i})$ (with \mathbf{I}_{n_i} an identity 214 matrix of dimension n_i = the number of repeated measurements in a subject) 215 and $\varepsilon_{(2)i} \sim N(\mathbf{0}, \tau^2 \mathbf{H}_i)$ (with \mathbf{H}_i the serial correlation matrix that only 216 depends on *i* through the number of repeated measurements n_i and the time 217 points j and k at which the measurements are taken). The (j, k) element h_{ijk} 218 of \mathbf{H}_i can then be modeled as $h_{ijk} = g(|t_{ij} - t_{ik}|)$ for a decreasing function g. 219 Two frequently used functions are the exponential and Gaussian correlation 220 functions, defined as $g(u_{j,k}) = \exp(-\phi u_{j,k})$ and $g(u_{j,k}) = \exp(-\phi u_{j,k}^2)$, 221 respectively. 222

223 4.2 Case study analysis

The mean structure of the model The LMMs that will be fitted to the case study dataset include an intercept, measurement moment, PEEP, and Cycle as fixed effects. PEEP and Cycle are dummy-coded with 5 and dummies, respectively. The relation between measurement point and the ZSV outcome is modeled as $\beta_1 t^3 + \beta_2 t^2 + \beta_3 t^2 \log(t)$ (see the Fractional polynomial section).

The covariance (correlation) structure of the model In the analyses below, three LMMs with the same fixed-effect structure (see previous paragraph) but different variance structures will be considered.

Model 1 is a random intercept model, i.e., a LMM that only contains a random intercept in the random part of the model:

$$Y_{ij} = \mu_{ij} + b_{0i} + \varepsilon_{ij},\tag{5}$$

where Y_{ij} is the observed endpoint at measurement time j for subject i, μ_{ij} is the mean as a function of the fixed effects, b_{0i} is the random intercept, and ε_{ij} is the residual. Based on this model, the reliability of the repeated observations taken at measurement times t_k and t_j can be estimated as (for details, see [12]):

$$R(t_j, t_k) = R = \frac{d}{d + \sigma^2},\tag{6}$$

where d is the variance of the random intercept and σ^2 is the residual variance. As can be seen in Eq. (6), the random intercept model assumes that any two observations measured at different times have the same R. This assumption is often not realistic when repeated measures are considered, i.e., measurements that are closer in time can be expected to be more strongly correlated than measurements that are more distant in time.

Therefore, Model 2 extends Model 1 by adding a serial correlation component:

$$Y_{ij} = \mu_{ij} + b_{0i} + \varepsilon_{(1)ij} + \varepsilon_{(2)ij}, \tag{7}$$

where μ_{ij} , b_{0i} are the same as in Model 1 and $\varepsilon_{(1)ij}$, $\varepsilon_{(2)ij}$ are measurement error and serial correlation components, respectively. Based on Model 2, the reliability of the repeated observations taken at measurement times t_k and t_j can be estimated as (for details, see [12]):

$$R(t_{j}, t_{k}) = R(u_{jk}) = \frac{d + \tau^{2} \exp\left(\frac{-u_{jk}^{2}}{\rho^{2}}\right)}{d + \tau^{2} + \sigma^{2}},$$
(8)

where $u_{jk} = t_k - t_j$, $\sigma^2 = \operatorname{Var}(\varepsilon_{(1)i})$ and $\tau^2 = \operatorname{Var}(\varepsilon_{(2)i})$. Model 2 thus no longer assumes that R remains constant for all pairs of measurements. Instead, it models R as a function of the time lag u_{jk} between two measurements. As can be seen, a stronger serial effect (ρ^2) leads to a faster decreasing $R(u_{jk})$.

Finally, Model 3 further extends Model 2 by including a random slope for measurement moment:

$$Y_{ij} = \mu_{ij} + b_{0i} + b_{1i}t_j + \varepsilon_{(1)ij} + \varepsilon_{(2)ij}, \qquad (9)$$

where μ_{ij} , b_{0i} , $\varepsilon_{(1)ij}$, $\varepsilon_{(2)ij}$ are the same as in Models 1 and 2, and b_{1i} is the random slope for measurement moment. Based on Model 3, the reliability of the repeated observations measured at times t_k and t_j can be estimated as (for details, see [12]):

$$R(t_j, t_k) = \frac{\mathbf{z}_j \mathbf{D} \mathbf{z}'_k + \tau^2 \exp\left(\frac{-u_{jk}^2}{\rho^2}\right)}{\sqrt{\mathbf{z}_j \mathbf{D} \mathbf{z}'_j + \tau^2 + \sigma^2} \sqrt{\mathbf{z}_k \mathbf{D} \mathbf{z}'_k + \tau^2 + \sigma^2}},$$
(10)

where $u_{jk} = t_k - t_j$, and \mathbf{z}_j , \mathbf{z}_k are the design rows in \mathbf{Z} corresponding to time j and k, respectively. As can be seen in Eq. (10), Model 3 no longer assumes that measurements taken at different time points but with the same time lag have the same R. Instead, it provides estimates of reliability for all pairs of measurements.

Table 2 summarizes the covariance structures that are used in the different models and their impact on the estimated R.

270

 $_{\rm 271}~$ » Table 2 about here «

272 4.2.1 Model 1: random intercept model

When Model 1 was fitted to the case study dataset, it was obtained that $\hat{d} =$ 1901.611 and $\hat{\sigma}^2 = 2413.022$, yielding $\hat{R} = 0.441$ (see Eq. (6)). A CI around \hat{R} can be computed by using a (non-parametric) bootstrap or the Delta method (for details, see the Web Appendix Part VI). The bootstrap-based 95% CI (using 500 bootstrap samples) equaled [0.198; 0.618]. The Delta methodbased CI was similar and largely overlapped, i.e., [0.189; 0.636]. Figure 4 (top left) illustrates the results (the bootstrap-based CI is shown).

Overall, it can be concluded that \widehat{R} is moderate and that there is substantial uncertainty in \widehat{R} (which is not surprising given the small number of ²⁸² animals in the study).

283

» Figure 4 about here «

²⁸⁵ 4.2.2 Model 2: random intercept and serial correlation

When Model 2 was fitted to the data of the case study, the estimated covariance parameters were $\hat{d} = 1349.650$, $\hat{\tau}^2 = 2489.351$, $\hat{\rho} = 3.581$, and $\hat{\sigma}^2 = 382.795$. Thus, after correction for the fixed effects, the covariance parameter estimates showed considerable remaining serial components.

Figure 4 (top right) shows the estimated $R(u_{ik})$ (see Eq. (8)) and their 290 95% CIs based on a bootstrap (the Delta method-based CIs were similar; data 291 are shown in the Web Appendix Part I). As can be seen, the estimated R were 292 high for small time lags (e.g., $\hat{R}(u_{ik} = 0) = 0.865$ and $\hat{R}(u_{ik} = 1) = 0.751$) 293 and subsequently decreased until they remained essentially constant at $R \approx$ 294 0.320 for measurements with time lags of about $u_{ik} = 10$ and higher. It can 295 also be observed that the CIs around $R(u_{ik})$ were narrower for measurements 296 with smaller time lags (e.g., for time lags $u_{ik} = 0$ and $u_{ik} = 1$, the $CI_{95\%} =$ 297 [0.817, 0.906] and $CI_{95\%} = [0.654, 0.836]$, respectively) and subsequently 298 widened until they remained stable around time lag $u_{ik} = 10$ with $CI_{95\%} =$ 299 [0.045, 0.530].300

³⁰¹ 4.2.3 Model 3: random intercept, slope, and serial correlation

When Model 3 was fitted to the data of the case study, the estimated covariance parameters were $\hat{\tau}^2 = 1952.970$, $\hat{\rho} = 3.290$, $\hat{\sigma}^2 = 373.043$, and

$$\widehat{\mathbf{D}} = \left(\begin{array}{cc} 3219.869 & -77.377 \\ -77.377 & 3.686 \end{array}\right)$$

As noted earlier, based on Model 3 the estimated $R(t_k, t_i)$ are different for 304 all pairs of measurements (see Eq. (10)). Figure 4 (bottom) shows the re-305 sults graphically. In this figure, the utmost left line (marked with t_1) depicts 306 the estimated $R(t_1, t_j)$, i.e., the estimated reliabilities of ZSV taken at mea-307 surement times 1 and 2–45. The line next to that one shows the estimated 308 $R(t_2, t_i)$, etc. The figure shows that $\widehat{R}(t_k, t_i)$ is high when the time lag u is 309 small and flattens out for longer time lags. Further, depending on the partic-310 ular pair of measurement moments (t_k, t_j) that is considered, the slope and 311 amount of decline in $R(t_k, t_i)$ as a function of time lag differs. For example, 312

when considering $\widehat{R}(t_1, t_j)$, it can be seen that the estimated reliabilities decline particularly strong for the first few subsequent measurements (say, until about t_8) and continue to decline for all t_j afterwards at a slower pace. Instead, for $\widehat{R}(t_{20}, t_j)$ there is only a substantial decline in the estimated reliabilities for the first few subsequent measurements (say, until about t_{25}) after which the estimated reliabilities remain essentially constant.

Based on Model 3, estimates of reliability are provided for each pair of measurements, and the same obviously holds for the CIs. To avoid cluttered figures, no CIs were added to Figure 4 (bottom). By means of illustration, Figure 5 provides 95% bootstrap-based CIs for $\hat{R}(t_1, t_j)$ (left) and $\hat{R}(t_{20}, t_j)$ (right). As can be seen, the CIs increase as a function of time and tend to be wider for $\hat{R}(t_{20}, t_j)$ than for $\hat{R}(t_1, t_j)$ (as expected).

325 326

» Figure 5 about here «

327 4.2.4 Selecting the most appropriate model

Based on the likelihood ratio (LR) test statistic G^2 , the fit of Models 1–3 328 can be formally compared (for details, see [7]). G^2 is equal to -2 times 329 the difference of the log likelihoods of the models being compared. Before 330 discussing the results for the case study, some general remarks are useful. 331 First, when interest is in testing the need for including random effects in the 332 model, the usual procedure where the test statistic G^2 is compared to a χ^2 333 distribution with the number of degrees of freedom equal to the difference 334 in the model parameters to be estimated is no longer valid. For example, 335 consider the situation where interest is in testing whether one or two random 336 effects are needed (Model 2 versus Model 3). This corresponds to testing 337 that $d_{12} = d_{21} = d_{22} = 0$. To test this hypothesis, a *mixture* with equal 338 weights 0.5 for χ_1^2 and χ_2^2 is needed (denoted by $\chi_{1:2}^2$), because the variance 339 d_{22} cannot be negative and thus the hypothesis test of interest is on the 340 boundary of the parameter space (for details, see [7]). Second, the results of 341 the LR tests should be interpreted with caution because of the small sample 342 size in the case study. Alternative testing procedures that are based on 343 permutation tests (see e.g., [13]) could provide a more viable alternative, but 344 these methods are beyond the scope of the present paper. Third, the valid 345 use of LR tests typically requires that the models are fitted using Maximum 346 Likelihood estimation. The results provided above used Restricted Maximum 347 Likelihood (REML), but valid LR tests for comparing nested models with 348

different covariance structures can still be obtained under REML estimation
when the models that are compared have the same mean structure [7] – which
was the case here, see above.

The log likelihood values for Models 1–3 are shown in Table 3. As can 352 be seen, the random intercept model with serial correlation (Model 2) fitted 353 the data significantly better than the random intercept model with no serial 354 correlation (Model 1), p < 0.001. This test thus rejects the null hypothesis 355 that there is no serial correlation process, i.e., it can be concluded that ob-356 servations that are closer in time are stronger correlated than observations 357 that are more distant in time. Further, adding a random slope to the random 358 intercept model with serial correlation (Model 3 versus Model 2) significantly 359 improves the model fit, p = 0.015 – though the gain was quite modest. 360

Model 3 is the model with the largest likelihood. It would be preferred 361 if we would solely rely on statistical arguments. However, from an applied 362 perspective - i.e., also considering the practical usefulness of the results for 363 a clinician or researcher – Model 2 is arguably to be preferred over Model 364 3 because the former leads to reliability estimates that only depend on the 365 time lag between two measurements. In contrast, Model 3 yields different 366 reliability estimates for all possible pairs of measurements. Model 2 thus 367 provides a much more parsimonious result compared to Model 3 – whilst the 368 fit of both models is roughly comparable. Notice that the likelihood ratio 369 tests identify the best fitting model among the models that were under con-370 sideration. However, when a model has been selected, the question remains 371 whether this model fits the data sufficiently well. Residuals and influence 372 diagnostics are useful in this respect. In Part VII of the Supplementary Ma-373 terials, a residual analysis is conducted and the extent to which particular 374 animals exert a strong influence on the results (i.e., the REML distances of 375 the models, the estimated fixed-effects parameters, the estimated covariance 376 components, and the estimated reliability coefficients) is evaluated. Overall, 377 the impact of excluding an animal on the results was relatively small for 378 Models 2–3. For Model 1, the impact of deleting an animal on the results 379 was more substantial. Further, the residual analysis showed that there were 380 no major departures of normality. 381

382

383 384

st Table 3 about here «

385 5 Discussion

The conventional methods to estimate reliability (e.g., the well-known Pear-386 son correlation coefficient) require assumptions that are often not met in 387 real-life studies (e.g., parallel measurements, equally spaced test-retest inter-388 vals, etc.). The main aim of the current paper was to present a general and 389 flexible approach to estimate reliability that is based on LMMs. It was shown 390 that this approach can be successfully applied even in a 'challenging' dataset 391 like in the presented case study – where the number of independent subjects 392 is low, different subjects have a different number of repeated observations, 393 and several covariates have to be taken into account. Overall, the analysis 394 of the case study suggested that the reliability of ZSV was high (and its CIs 395 narrow) when the time lag was small. For larger time lags, the reliability 396 estimates decreased and their CIs widened. 397

Some critical remarks are in place. First, despite the major differences be-398 tween the conventional and the LMM-based methods to estimate reliability, 399 there are also some obvious similarities. For example, the expressions to esti-400 mate reliability based on Model 1 (see Eq. (6)) and the conventional approach 401 (see Eq. (2)) are very similar (i.e., both are ratios of variances). However, 402 a fundamental difference between both methods is that the LMM-based ap-403 proach does not require the parallel measurement assumption. The reason 404 for this is that the mean and variance structures can be clearly separated in 405 LMMs (see above). For example, when the means at different time points 406 are different (as was observed in the case study, see Figure 1), systematic 407 effects of time and other covariates can be taken into account by including 408 them into the fixed-effect part of the model (as was done here). In essence, 409 the main difference between the conventional and LMM-based approaches 410 to estimate reliability is that the former requires a set of assumptions that 411 are taken care of in the study design, whereas the latter takes care of these 412 assumptions through modelling at the analysis stage [3]. There is however a 413 price to pay for the increased flexibility of the LMM-based approach, i.e., it 414 requires substantially more complex statistical analyses compared to the con-415 ventional methods to estimate reliability. We tried to circumvent this issue 416 by developing an R package (*CorrMixed*) that allows for obtaining reliability 417 estimates based on Models 1–3 in a relatively straightforward way. The Web 418 Appendix (Parts IV and V) provides full details on how the analyses can be 419 conducted in practice. 420

421 Second, in the present paper the focus was entirely on the random effect

structure of the models because we were interested in estimating the reliabil-422 ity of the outcomes. Apart from estimating reliability, medical practitioners 423 are also often interested in obtaining so-called normative data. Normative 424 data are used to convert a patient's 'raw' outcomes into relative measures 425 that reflect the proportion of demographically-matched healthy controls in 426 the population who have a lower outcome value compared to this patient. 427 A well-known example are growth curves of young children. Such normative 428 data (nomograms) for repeated measurements can be obtained without any 429 substantial additional effort using the same type of models that were fitted 430 in the present paper. The only difference is that the focus will then be on the 431 fixed-effect part of the model – rather than on the random effect structure 432 (for details, see [14]). 433

Third, the outcome that was considered in the case study was a normally distributed (Gaussian) variable. One may also be interested in estimating the reliability of repeated measurements of outcomes of a different distributional nature, e.g., binary (yes/no, health/sick) or categorical ordered outcomes. Such extensions are possible, but not trivial. The interested reader is referred to Vangeneugden *et al.* [15].

Fourth, in the analysis of the case study, the fixed-effect structures were 440 kept constant for Models 1 to 3 (because we were primarily interested in 441 evaluating the impact of different random-effect structures on the estimated 442 reliabilities). In the Web Appendix (Part III), a sensitivity analysis is con-443 ducted where the impact of using different plausible fixed-effect structures 444 on the estimated reliabilities is evaluated. Overall, the analyses indicated 445 that the estimated reliabilities are not sensitive to the fixed-effect part of the 446 model (provided that the mean structure of the model is supported by the 447 data). 448

Finally, in the present paper no time-varying covariates (other than mea-449 surement occasion itself) were considered, but depending on the study at 450 hand it may be useful to include such covariates. For example, consider a 451 setting where one is interested in estimating the reliability of a psychiatric 452 rating scale that was scored by different physicians at the different mea-453 surement moments. When only a limited number of raters are involved in 454 the study, the methodology that was proposed above can still be used in a 455 straightforward way. Indeed, one can then simply include rater as a (dummy-456 coded) fixed-effect in the mean structure of the model. On the other hand, 457 when the number of raters is large, it is more sensible to include rater in the 458 random-effect part of the model. Such a model cannot be fitted in the cur-459

rent version of the *CorrMixed* package, but it is straightforward to fit such a model using SAS.

On a related note, in the present paper interest was primarily in the es-462 timation of the reliability of a single outcome that was repeatedly measured 463 within the same subject. It might also be of interest to estimate how strongly 464 the vectors of two outcomes are correlated with each other. For example, 465 consider a setting where two raters assess all patients at all measurement 466 moments. Here, it would be natural to study the correlation between the 467 vectors of scores to evaluate the level of agreement between the two raters. 468 Or, as another example, consider a setting where there are two alternative 469 measurement procedures for the same latent variable. When one of the two 470 measurement procedures is more 'difficult' to conduct (e.g., is more expen-471 sive, more painful for the patient, requires more time to obtain the test 472 results, etc), it may be of interest to estimate the correlation between the 473 measurements obtained by both procedures. Indeed, when it can be shown 474 that there is a high correlation between the vectors of outcomes, the 'easier' 475 measurement procedure may replace the more difficult one - in the same 476 spirit as is done when a surrogate endpoint is used to replace the true end-477 point in a clinical trial (individual-level surrogacy; for details see [16]). The 478 quantification of the correlation between two vectors of outcomes is however 479 beyond the scope of the present paper, as different statistical techniques are 480 needed to estimate this quantity (see e.g., [17]). 481

Acknowedgements

Financial support from the IAP research network #P7/06 of the Belgian Government (Belgian Science Policy) is gratefully acknowledged. This project has received funding from the European Union's 7th Framework Programme for research, technological development and demonstration under the IDEAL Grant Agreement no 602552.

Web Appendix

A Web Appendix is available that contains (i) the Delta method-based 95% CIs of the estimated reliability coefficients for ZSV, (ii) a sensitivity analysis where the impact of 2 clinically deviating animals on the results is examined,

(iii) a sensitivity analysis where the impact of using a different fixed-effect structure on the results is examined, (iv) details on how the newly developed R package *CorrMixed* can be used to estimate reliability, (v) details on how reliability can be estimated using SAS, (vi) details on the computation of the Delta method-based CIs for \hat{R} , and, (vii) the results of a residual analysis.

References

- Lord FM, Novick MR. Statistical theories of mental test scores. Addison-Welsley Publishing Company, Reading, MA; 1968.
- [2] Spearman C. The proof and measurement of association between two things. The American Journal of Psychology 1904; 15:72-101.
- [3] Laenen A. Psychometric Validation of Continuous Rating Scales from Complex Data; 2008. Unpublished PhD thesis. Retrieved from http://ibiostat.be/publications/phd/annouschkalaenen.pdf
- [4] Molenberghs G, Kenward M. Missing data in clinical studies. New York: John Wiley & Sons; 2007.
- [5] Rubin DB. Inference and missing data. *Biometrika* 1976; **63**:581-592.
- [6] Fleiss JL. Design and analysis of clinical experiments. Wiley: New York; 1986.
- [7] Verbeke G, Molenberghs G. Linear Mixed Models for Longitudinal Data. New York: Springer-Verlag; 2000.
- [8] Pikkemaat R, Lundin S, Stenqvist O, Hilgers, RD, Leonhardt, S. Recent advances in and limitations of cardiac output monitoring by means of electrical impedance tomography. *Anesthesia & Analgesia* 2014; 119:76-83.
- [9] Royston P, Altman, DG. (1994). Regression using fractional polynomials of continuous covariates: parsimonious parametric modelling (with discussion). Journal of the Royal Statistical Society. Series C (Applied Statistics) 1994; 43:429-467.
- [10] Glaser D, Hastings RH. (2011). An introduction to multilevel modeling for anesthesiologists. Anaesthesia & Analgesia 2011; 113:877-887.
- [11] West BT, Welch KB, Galecki AT. Linear Mixed Models. A practical guide using statistical software (2nd Ed.). New York: CRC Press, Taylor & Francis Group; 2015.

- [12] Vangeneugden T, Laenen A, Geys H, Renard D, Molenberghs G. Applying linear mixed models to estimate reliability in clinical trial data with repeated measurements. *Controlled Clinical Trials* 2004; 25:13-30.
- [13] Lee OE, Braun T. Permutation tests for random effects in linear mixed models. *Biometrics* 2012; 68:486-493.
- [14] Van der Elst W, Molenberghs G, Van Boxtel MPJ, Jolles J. Establishing normative data for repeated cognitive assessment: a comparison of different statistical methods. *Behavior Research Methods* 2013; 45:1073-1086.
- [15] Vangeneugden T, Molenberghs G, Laenen A, Geys H, Beunckens C, Sotto C. Marginal Correlation in Longitudinal Binary Data Based on Generalized Linear Mixed Models. *Communications in Statistics. Theory and Methods* 2010; **39**:3540-3557.
- [16] Burzykowski T, Molenberghs G, Buyse M. The Evaluation of Surrogate Endpoints. New York: Springer-Verlag; 2005.
- [17] Roy A. Estimating correlation coefficient between two variables with repeated observations using mixed effects model. *Biometrical Journal* 2006; 48:286-301.

Tables

	· riaconomar port	10111101 1 00 01
M	power p_m	AIC
1	-0.5	3788.703
2	0.5, 0.5	3786.096
3	2, 2, 3	3775.281
4	0.5,1,2,2	3776.389
5	-2, -2, 0, 2, 3	3778.221

Table 1: Fractional polynomial results.

Table 2: Summary of the covariance structures used in Models 1–3, and the impact on the estimated reliabilities.

Model	Estimated reliabilities R				
Model 1: Random Intercept	\widehat{R} is identical for all pairs (t_j, t_k)				
Model 2: Random intercept and serial component	\widehat{R} only depends on the time lag $u_{jk} = t_k - t_j$				
Model 3: Random intercept, slope, and serial component	\widehat{R} is different for all pairs (t_j, t_k)				
Note. t_j = measurement at time j .					

Table 3: Fit indices of the different models for the ZSV outcome.								
	# Pars.		$\log L$	G^2	Test	p		
	Rand.	Ser.						
Model 1	1	0	-2328.910					
Model 2	1	2	-2125.135	407.551	Model 2 vs. 1: χ^2_2	< 0.001		
Model 3	3	2	-2121.399	7.472	Model 3 vs. 2: $\chi^2_{1:2}$	0.015		
Note. $\log L = \log$ likelihood, $G^2 = -2$ the difference of two log likelihood								

values. Rand. = random effect parameters, ser. = serial components.

Figures



Figure 1: Individual profiles (grey lines) and mean values (black line) of the ZSV outcome as a function of time of measurement.



Figure 2: Number of observations for the ZSV outcome as a function of time of measurement.



Figure 3: Observed means as a function of time of measurement (solid line) and fitted fractional polynomial of degree m = 3 (dashed line).



provided to avoid a cluttered figure. The utmost left line marked with t_1 depicts the estimated correlations Model 1 (upper left), Model 2 (upper right) and Model 3 (bottom). For Model 3, no Confidence Intervals are Figure 4: Estimated reliabilities (solid lines) and 95% Confidence Intervals (dashed lines) for ZSV based on between t_1 and all other measurements, the line next to that one depicts the correlations between t_2 and measurements 2-45, and so on.



Figure 5: $\widehat{R}(t_1, t_j)$ (left) and $\widehat{R}(t_{20}, t_j)$ (right) based on Model 3 and their 95% Confidence Intervals for the ZSV outcome.