

Generalizability in non-Gaussian longitudinal clinical trial data based on generalized linear mixed models

Peer-reviewed author version

VANGENEUGDEN, Tony; MOLENBERGHS, Geert; LAENEN, Annouschka;  
ALONSO ABAD, Ariel & GEYS, Helena (2007) Generalizability in non-Gaussian longitudinal clinical trial data based on generalized linear mixed models. In: JOURNAL OF BIOPHARMACEUTICAL STATISTICS, 18(4). p. 691-712.

DOI: 10.1080/10543400802071386

Handle: <http://hdl.handle.net/1942/9221>

# Generalizability in Non-Gaussian Longitudinal Clinical Trial Data Based on Generalized Linear Mixed Models

Tony Vangeneugden<sup>1,2</sup>      Geert Molenberghs<sup>2</sup>  
Annouschka Laenen<sup>2</sup>      Ariel Alonso<sup>2</sup>  
Helena Geys<sup>2,3</sup>

<sup>1</sup> Tibotec, Johnson & Johnson, Mechelen, Belgium

Email: tvangene@tibbe.jnj.com

<sup>2</sup> Hasselt University, Center for Statistics, Diepenbeek, Belgium

<sup>3</sup> Janssen Pharmaceutica, Johnson & Johnson, Beerse, Belgium

## Abstract

This work investigates how generalizability, an extension of reliability, can be defined and estimated based on longitudinal data sequences resulting from, for example, clinical studies. Useful and intuitive approximate expressions are derived based on generalized linear mixed models. Data from four double-blind randomized clinical trials in schizophrenia motivate the research and are used to estimate generalizability for a binary response parameter.

**Keywords:** *Binary data; Intraclass correlation; Generalizability; Random effects; Reliability; Variance components.*

## 1 Introduction

Many measurements in clinical research are based on clinicians' observations, are therefore prone to error, and hence call for assessment of observer reliability and agreement. The latter terms are often used interchangeably but, in principle, should be considered different concepts. *Reliability* coefficients express the ability to differentiate among subjects and are ratios of variances; in classical terms, the variance attributed to the difference among subjects divided by the total variance (Shrout and Fleiss 1979). *Agreement* refers to conformity, with corresponding parameters determining whether the same value is achieved if a measurement were performed twice,

either by the same or different observers. In homogeneous populations, one can imagine that reliability be low while agreement be high; in a heterogeneous population, reliability and agreement measures allegedly will correspond well (Stratford 1989). The parameters for assessment of observer reliability and agreement differ according to the scale of measurement. For nominal and ordinal categorical measurements, the  $\kappa$ -coefficient and the weighted  $\kappa$ -coefficient, respectively, are measures of agreement. In case of continuous data, the intraclass correlation coefficient (ICC) is commonly used to measure observer reliability, although ICC-type quantities can be defined for binary and ordinal data as well (Fleiss 1981).

As stated by Fleiss (1986): “The most elegant design of a clinical study will not overcome the damage by unreliable or imprecise measurement.” In clinical trials, one typically wants to differentiate among treatments. If reliability is low, the ability to differentiate between the different subjects in the different treatment arms decreases. Fleiss describes as consequences of *unreliability*: attenuation of correlation in studies designed to estimate correlation between variables with poor reliability, biased sample selection in clinical studies where patients are selected based on attaining a minimum level of an unreliable measurement, and, last but not least, an increased sample size for trials with a primary, low-reliability parameter. For the latter, one can easily show that, for a paired  $t$ -test, the required sample size becomes  $n = n^*/R$  where  $R$  denotes the reliability coefficient and  $n^*$  is the required sample size for the true score, i.e., the required sample size when responses are measured without error. It is thus clear that a high reliability is important to the clinical trialist.

When the biostatistician and clinician are designing a new clinical study, they should have good information on the reliability of the planned measurements. Most often, the strategy is to use a scale that has been validated before and for which intra-rater (i.e., test-retest), inter-rater reliability, and internal consistency were established. The validation is usually done on a selected small sample from the population for which the scale is intended. If the population of the trial is different, a new battery of reliability and validity testing might be warranted. Now, the classical framework may be deficient for the clinical-trial setting since conventionally an observation is

assumed to be a combination of an individual's *true* score and random measurement error. The assumption that all variance in scores can be divided into true and error variance may come across as a little simplistic. At the same time, once the trial is finished and reported, it is astonishing how little attention is given to the observed reliability of a given scale. The focus usually is on estimating treatment effects and their significance. As a result, rarely is there any reflection on how reliable the scale was or how large the observed measurement error.

Vangeneugden *et al* (2004) proposed a framework for studying *trial- or population-specific reliability* using clinical-trial data. The appeal of this extension notwithstanding, next to the true score of an individual, multiple potential sources of error can exist. The goal is then, of course next to the main goal to study treatment differences, to also obtain the most precise estimate of the score a person would have if there were no sources of error contaminating the results. Each one among the variety of forms reliability can take, such as inter-rater reliability, test-retest reliability, and internal consistency, identifies and quantifies only one source of error variance at a time. *Generalizability theory* (GT, Cronbach, Rajaratnam, and Gleser 1963) enables considering all sources of variability simultaneously. The essence of the theory is the recognition that, in any measurement situation, there are multiple sources of error variance. The goal is to try and identify, measure, and thereby possibly find strategies to reduce the influence of these sources on the measurement under investigation (Shavelson, Webb, and Rowley 1989). Imagine that we could identify the most likely sources of error in a measurement of some characteristic of a person. We then have defined our “universe” (Cronbach *et al* 1972, Brennan 1992) of possible observations. If we subsequently proceed to average each person's score over all of these possible conditions, an unbiased estimate would result of that person's score over the universe as we have defined it. Note that there is no pretense that this is the “true” score; rather, it is conditional on the universe considered. Of course, more than one choice of universe can be considered.

In the context of clinical trials, by investigating sources of error, such as, for instance, country or sub-category of diagnosis, the clinical trialist could learn a lot about performance of scales or other measurements in certain subgroups and what the impact of such factors is on reliability.

Vangeneugden *et al* (2005) applied such generalizability concepts to interval-scaled data from clinical trials, for which it is natural to assume a Gaussian distribution.

The present work extends generalizability to non-Gaussian outcomes; in particular, our focus will be on binary data. Even in the univariate case, there are fundamental differences between Gaussian and non-Gaussian outcomes, since the latter usually require non-linear models, also exhibiting important differences in the relationship between mean and variance. Furthermore, repeated binary data are frequently encountered in clinical trials but pose challenges for model formulation. One distinguishes between marginal and random-effects model families and, unlike in the Gaussian situation, there is no easy relationship between both. An example of the marginal family is generalized estimating equations (GEE, Liang and Zeger 1986), whereas the generalized linear mixed model (GLMM, Breslow and Clayton 1993) is likely the most prominent random-effects model (Molenberghs and Verbeke 2005). Whereas GEE is convenient and frequently used, it models the marginal regression function, treating the second and higher-order moments as nuisance parameters, which limits its use when the correlation is of scientific interest, e.g., in view of the ICC. The GLMM, on the other hand, has a full likelihood basis, but fails to produce the marginal correlations in an easy fashion, owing to the presence of a non-linear link function, combined with a non-trivial mean-variance relationship, forcing the variance to change with the mean and hence with the regressors (Molenberghs and Verbeke 2005, Chapter 16). In spite of these considerations, we will show the GLMM provides a viable framework when correlations are of interest, with particular emphasis on the use of generalizability theory.

In Section 2, the motivating case study is introduced, while methodology is described in Section 3. In Section 4, we will estimate reliability and generalizability of a binary response variable, thereby underscoring the versatile use that can be given to the generalizability framework.

## 2 Motivating Studies

Consider individual patient data from four double-blind randomized clinical trials, comparing the effects of risperidone to conventional anti-psychotic agents for the treatment of chronic schizophrenia. Schizophrenia has long been recognized as a heterogeneous disorder with patients suffering from both “negative” and “positive” symptoms. Negative symptoms are characterized by deficits in social functions such as poverty of speech, apathy and emotional withdrawal. Positive symptoms entail more florid symptoms such as delusions and hallucinations, which are superimposed on the mental status. Several measures can be considered to assess a patient’s global condition. The *Positive and Negative Syndrome Scale* (PANSS) consists of 30 items that provide an operationalized, drug-sensitive instrument, which is highly useful for both typological and dimensional assessment of schizophrenia (Kay, Fiszbein, and Opler 1987). Classical reliability of the PANSS has been studied previously (Kay, Opler, and Lindenmayer 1988, Bell *et al* 1992, Peralta and Cuesta 1994). It will be used in Section 4.5. The *Clinical Global Impression* (CGI) of overall change versus baseline is a 7-grade scale used by the treating physician to characterize how well a subject has improved since baseline. The levels are: ‘very much improved,’ ‘much improved,’ ‘minimally improved,’ ‘no change,’ ‘minimally worse,’ ‘much worse,’ ‘very much worse.’ Clinical response is often defined as a CGI score of ‘very much improved or ‘much improved.’ Since the label in most countries recommend doses ranging within 4–6 mg/day, we include in our analysis only patients who received either these doses of risperidone or an active control (haloperidol, perphenazine, or zuclopenthixol). Depending on the trial, treatment was administered for a duration of 6–8 weeks. For example, in the international trials by Peuskens *et al* (1995), Marder and Meibach (1994), and Hoyberg *et al* (1993) patients received treatment for 8 weeks, while in the study by Huttunen *et al* (1995) patients were treated over a period of 6 weeks. The sample sizes were 453, 176, 74, and 71, respectively. Measurements were taken at weeks 1, 2, 4, 6, and 8. In Section 4, the pooled data from these trials will be analyzed; this decision is defensible given the trials are compromised of similar patients.

### 3 Methodology

We will first present a general outline of the concepts of reliability and generalizability from a classical view-point, i.e., when the outcomes are assumed normally distributed (Sections 3.1 and 3.2). Thereafter, we will sketch the generalized linear mixed model paradigm, which will enable us to deal with binary and other non-Gaussian outcomes (Section 3.3). This then offers a framework within which reliability and generalizability can be derived based on non-Gaussian longitudinal data from clinical trials or other studies, not specifically designed in view of reliability or generalizability (Section 3.4).

To fix ideas, let us give an example as to how the observed clinical trial data are typically decomposed:

$$Y_{pdt} = h(\mu + b_p + \mu_d^D + \mu_t^T + \mu_{dt}^{DT}) + \varepsilon_{pdt}, \quad (1)$$

where  $h(\cdot)$  is a known link function. Further,  $b_p$  denotes the random effect for patient  $p = 1, \dots, N$ ,  $\mu_d^D$  the fixed time effect at day  $d = 1, \dots, n_p$ ,  $\mu_t^T$  the fixed effect of treatment  $t = 1, \dots, T$ ,  $\mu_{dt}^{DT}$  their interaction. Finally,  $\varepsilon_{pdt}$  refers to the residual error, the distribution of which is chosen in accordance with the outcome type. For example, when  $Y_{pdt}$  is a binary indicator, it is customary to adopt for  $h(\cdot)$  the antilogit function and for  $\varepsilon_{pdt}$  the Bernoulli distribution with success probability  $h(\mu + b_p + \mu_d^D + \mu_t^T + \mu_{dt}^{DT})$ . When other design levels are present, e.g., country or center, Model (1) can be extended in a straightforward fashion and various instances will be given in subsequent sections.

#### 3.1 Reliability

In classical test theory, reliability frequently materializes as the intraclass correlation coefficient (ICC). For instance, if one wishes to estimate test-retest reliability in case of Gaussian data, the outcome of a test can be modeled as

$$Y_{pd} = \mu + b_p + \mu_d^D + \varepsilon_{pd}, \quad (2)$$

where  $\mu$  is an overall intercept,  $b_p \sim N(0, \sigma_p^2)$  a random effect for patient,  $\mu_d^D$  a fixed effect for day of measurement, and  $\varepsilon_{pd} \sim N(0, \sigma_E^2)$  the corresponding measurement error. Then, the reliability is a function of two sources of variability, deriving from the patient and residual levels, respectively:

$$\hat{R} = \frac{\hat{\sigma}_p^2}{\hat{\sigma}_p^2 + \hat{\sigma}_E^2}, \quad (3)$$

This expression ought to be time-independent, because the classical theory assumes strictly parallel measurements (Shavelson, Webb, and Rowley 1989, p. 924). Of course, in what follows we will switch to generalizability, where this assumption and hence the corresponding restrictions, will no longer be made.

It is possible to show that  $R$  is the correlation between measurements of the same patient, on different but given days, i.e., conditioning on days and thereby keeping them fixed:

$$R = \text{Corr}(Y_{pd}, Y_{pd'} \mid d, d'), \quad (4)$$

with notation as in (2). For parallel measurements, this correlation coefficient indeed coincides with the ICC of reliability.

### 3.2 Generalizability

Generalizability theory (Cronbach *et al* 1972, Shavelson, Webb, and Rowley 1989, Brennan 1992, Streiner and Norman 1995) recognizes that, in virtually all measurement situations, there are multiple sources of error variance. The goal is to try and identify, measure, and thereby possibly find strategies to reduce the influence of these sources on the measurement in question. For instance, if one measures a patient, not only on different days but also by different raters, one could investigate both test-retest reliability and inter-rater reliability, assuming that rater and day of observation are the most important sources of error, in addition to residual measurement error. Classically, these two reliabilities result from different calculations and, when done in a model-based fashion, they result from different models. In this paper, our approach is to specify general, encompassing models, which contain specific choices, such as test-retest reliability or



inter-rater reliability, as special cases. This is useful for the linear-models case but perhaps even more so for the non-linear case, discussed in Section 3.3 and later, where the generalized linear mixed model framework will be used as the modeling basis for our developments. Indeed, in such as case, more complex correlation computations have to be undertaken. Embedding all conventional coefficients within a single framework then implies that the calculations need to be done only once. This also has strong advantages in terms of software implementation.

For now returning to Gaussian data, a linear version of (2), allowing for rater effects and treating day as random, is:

$$Y_{prd} = \mu + b_p^P + b_r^R + b_d^D + \varepsilon_{prd}, \quad (5)$$

where, in addition to effects already included,  $b_r^R$  now represents the random effect pertaining to rater  $r = 1, \dots, R$ . The associated sources of variability are denoted by  $\sigma_p^2$  for the patient level,  $\sigma_r^2$  for the rater level, etc. Model (5) enables estimation of the variances' magnitude stemming from the various sources of error, which are patient, rater, day, and residual in the example above. If the sources that we have identified are trivial, while we have missed any important source of error, then the residual variance will typically be large.

In GT terminology, 'person' is a so-called *facet of differentiation*, while 'rater' and 'day' are called *facets of generalization*. Note that the term *facet of differentiation* originates from Streiner and Norman (1995), whereas the term *object of measurement* is very commonly used. In what follows, the latter term will be employed. The levels of the facets of generalization are named *conditions*. It is common to use ANOVA for estimating the various variance components, which in turn lead to a *generalizability coefficient*, analogous to a reliability coefficient, found as the ratio of the estimated person variance component and a so-called estimated observed score variance.

GT distinguishes between decisions based on the relative standing of individuals and decisions based on the absolute value of a score (Shavelson, Webb, and Rowley 1989). Let us explain these in turn.

Error in *relative decisions* arises from all nonzero variance components associated with rank

ordering of individuals, other than the component for the object of measurement (persons). These can stem from effects present in Model (5) or from an extension of the model arising, for example, from including interaction terms:

$$Y_{prd} = \mu + b_p^P + b_r^R + b_d^D + b_{pr}^{PR} + b_{pd}^{PD} + b_{rd}^{RD} + \varepsilon_{prd}, \quad (6)$$

In such a model, variance components associated with the interaction of person with each facet or combinations of facets define sources of error. For Model (6), we additionally distinguish between  $\sigma_{pr}^2$ ,  $\sigma_{pd}^2$ , and  $\sigma_{prd}^2 = \sigma_e^2$ . So, if one wishes to generalize from a rating by one rater on a particular day to a rating by a different rater at another point in time, the following generalizability coefficient can be constructed as the ratio of the universe-score variance to the expected rater-score variance:

$$E_{\rho^2 \text{Rel}} = \text{Corr}(Y_{prd}, Y_{pr'd'} \mid r, r', d, d') = \frac{\sigma_p^2}{\sigma_p^2 + \sigma_{\text{Rel. Error}}^2} = \frac{\sigma_p^2}{\sigma_p^2 + \sigma_{pr}^2 + \sigma_{pd}^2 + \sigma_{prd}^2}, \quad (7)$$

having the form of an ICC. Indeed, it is possible to show that (7) can be derived as a conditional correlation coefficient where we condition on rater and day, while at the same time allowing each one of them to take on different values. Alternatively, we can derive a test-retest or an inter-rater reliability coefficient, by either merely generalizing over day of observation and fixing rater or by merely generalizing over rater and fixing day of observation:

$$R_{\text{test-retest, Rel}} = \text{Corr}(Y_{prd}, Y_{prd'} \mid r, d, d') = \frac{\sigma_p^2 + \sigma_{pr}^2}{\sigma_p^2 + \sigma_{pr}^2 + \sigma_{pd}^2 + \sigma_{prd}^2}, \quad (8)$$

$$R_{\text{inter-rater, Rel}} = \text{Corr}(Y_{prd}, Y_{pr'd} \mid r, r', d) = \frac{\sigma_p^2 + \sigma_{pd}^2}{\sigma_p^2 + \sigma_{pr}^2 + \sigma_{pd}^2 + \sigma_{prd}^2}. \quad (9)$$

Decisions based on the level of observed score, disregarding the performance of others, are called *absolute decisions*. All variance components associated with such a score, except the component for the object of measurement, are considered error. Then, (7) transforms to

$$\begin{aligned} E_{\rho^2 \text{Abs}} = \text{Corr}(Y_{prd}, Y_{pr'd'}) &= \frac{\sigma_p^2}{\sigma_p^2 + \sigma_{\text{Abs. Error}}^2} \\ &= \frac{\sigma_p^2}{\sigma_p^2 + \sigma_r^2 + \sigma_d^2 + \sigma_{pr}^2 + \sigma_{pd}^2 + \sigma_{rd}^2 + \sigma_{prd}^2}. \end{aligned} \quad (10)$$

Also here, (10) can be considered an ICC, this time conditioned neither on rater nor on day. Similar to the above, we can derive an *absolute* test-retest or inter-rater reliability coefficient:

$$R_{\text{test-retest, Abs}} = \text{Corr}(Y_{prd}, Y_{prd'}) = \frac{\sigma_p^2 + \sigma_r^2 + \sigma_{pr}^2}{\sigma_p^2 + \sigma_r^2 + \sigma_d^2 + \sigma_{pr}^2 + \sigma_{pd}^2 + \sigma_{rd}^2 + \sigma_{prd}^2}, \quad (11)$$

$$R_{\text{inter-rater, Abs}} = \text{Corr}(Y_{prd}, Y_{pr'd}) = \frac{\sigma_p^2 + \sigma_d^2 + \sigma_{pd}^2}{\sigma_p^2 + \sigma_r^2 + \sigma_d^2 + \sigma_{pr}^2 + \sigma_{pd}^2 + \sigma_{rd}^2 + \sigma_{prd}^2}. \quad (12)$$

This example, aimed at enhancing insight into the various uses of GT, was based on a simple so-called *crossed* design with two factors or facets of generalization,  $r$  for ‘rater’ and  $d$  for ‘day,’ of course in addition to the replication over ‘patient,’ the object of measurement or facet of differentiation, each one occurring at all levels of the other. GT can be used with more complex designs as well, for example, including more factors, and even in *nested* designs, exhibiting a more complex factor structure. As discussed in Streiner and Norman (1995), the general principle remains untouched: one begins by isolating the various sources of variance in the scores, and then generating a family of coefficients that depend on the particular factors and that are allowed either to vary or to remain fixed.

A first type of study, designed to estimate variance components underlying a measurement process, is called a *G-study*. Second, having generated the variance estimates, one can then study the impact on generalizability of such a decision as changing the number of observations or adding a further rater. Since this second type of study explores the impact of certain decisions, they are termed *Decision studies* or *D-studies*. Interestingly, D-studies can be undertaken solely using paper and pencil, or a computer. In planning a D-study, the decision maker defines the universe of generalization and specifies the proposed interpretation of the measurement. The goal is to identify important sources of variability in a particular measurement situation from the outset, and then one quantifies these sources. Ample detail, as well as insightful examples and perspectives on estimation, can be found in Shavelson, Webb, and Rowley (1989, in particular pp. 925–926).

Obviously, GT is broad and versatile. In the next section, we will show how this can be expanded by embedding it in the flexible generalized linear mixed model framework. Apart from dealing

with non-Gaussian outcomes, it will be possible to include further sources of variability, such as serial (temporal) correlation, commonly encountered in longitudinal studies, superimposed on the random-effects structure.

The fact that sample size does not intervene in many of the above formulas is intriguing and potentially confusing, especially for the reader used to conventional formulas that involve sample sizes. This is a consequence of using a flexible model set up, rather than confining attention to classical averaging, and in line with flexible models that can be found in Shavelson, Webb, and Rowley (1989) and Streiner and Norman (1995, their Section 9, Examples 1–3). The apparent absence of a sample size merely means it is implicitly present, as can be seen through the following analogy stemming from conventional linear regression. Indeed, the general linear regression model can be written as  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ , with  $\mathbf{Y}$  a vector of responses on a set of  $N$  individuals,  $\mathbf{X}$  an  $N \times p$  design matrix,  $\boldsymbol{\beta}$  a vector of unknown regression coefficients, and  $\boldsymbol{\varepsilon}$  an  $N \times 1$  vector of independent error terms. A general expression for the ordinary least squares or maximum likelihood estimator for  $\boldsymbol{\beta}$  is  $\hat{\boldsymbol{\beta}} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y}$ , which appears to be free of the sample size. Now, reducing the model to one with a random intercept only,  $\mathbf{X}$  becomes an  $N \times 1$  vector of ones, and  $\boldsymbol{\beta} = \mu$ , the overall mean. It then follows that:

$$\hat{\boldsymbol{\beta}} = \left[ (1 \dots 1) \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} \right]^{-1} (1 \dots 1) \begin{pmatrix} y_1 \\ \vdots \\ y_N \end{pmatrix} = \frac{1}{N} \sum_{i=1}^N y_i,$$

making the sample size reappear. Evidently, explicit expressions enhance clarity, but they sometimes limit generality. In sections to follow, we will switch to non-Gaussian data, in which case there usually are no analogous explicit expression. The advantage of a generic framework is then clear.

Another advantage of a general modeling framework is that it unifies what are classically perceived as different settings. For example, test-retest reliability (8) and inter-rater agreement (9) follow as special cases from the same setting, while at the same time other situations are possible, for which no classical counterparts exist, such as (7), as pointed out by Streiner and Norman (1995,

p. 135). This should not be viewed as a disadvantage but rather as an opportunity. Evidently, not every extension of this type will be relevant for a particular analysis, but the user can decide which particular version of the model to consider and, consequently, which coefficient to derive.

### 3.3 Generalized Linear Mixed Models

The generalized linear mixed model (GLMM, Breslow and Clayton 1993) has been the most frequently used random-effects model for non-Gaussian outcomes, although alternative paradigms, such as laid out in Lee, Nelder, and Pawitan (2006), exist and are of interest, too.

With notation similar to the one used in previous sections, let  $Y_{pd}$  be the outcome recorded on day  $d = 1, \dots, n_p$  for subject  $p = 1, \dots, N$ , and let  $\mathbf{Y}_p$  be the  $n_p$ -dimensional vector of all measurements available for subject (cluster)  $p$ . The GLMM assumes that, conditional on a  $q$ -dimensional random  $\mathbf{b}_p$ , allegedly drawn independently from a  $N(\mathbf{0}, H)$ , the outcomes  $Y_{pd}$  are independent with densities of the form

$$f_p(y_{pd}|\mathbf{b}_p, \boldsymbol{\beta}, \phi) = \exp \left[ \frac{y_{pd}\theta_{pd} - \psi(\theta_{pd})}{\phi} + c(y_{pd}, \phi) \right], \quad (13)$$

where the mean  $\mu_{pd}^{PD} = \partial\psi(\theta_{pd})/\partial\theta_{pd}$  is modeled through a linear predictor containing fixed regression parameters  $\boldsymbol{\beta}$  as well as subject-specific parameters  $\mathbf{b}_p$ , i.e.,  $h^{-1}(\mu_{pd}^{PD}) = h^{-1}[E(Y_{pd}|\mathbf{b}_p)] = \mathbf{x}'_{pd}\boldsymbol{\beta} + \mathbf{z}'_{pd}\mathbf{b}_p$  for a known link function  $h(\cdot)$ , with  $\mathbf{x}_{pd}$  and  $\mathbf{z}_{pd}$   $r$ -dimensional and  $q$ -dimensional vectors of known covariate values, respectively, with  $\boldsymbol{\beta}$  an  $r$ -dimensional vector of unknown fixed regression coefficients, and with  $\phi$  a scale parameter. Employing a natural link function (McCullagh and Nelder 1989), this becomes  $\theta_{pd} = \mathbf{x}'_{pd}\boldsymbol{\beta} + \mathbf{z}'_{pd}\mathbf{b}_p$ . Estimation of model parameters is slightly cumbersome, since no explicit formula for a subject's likelihood contribution exists. Therefore, one has to resort to numerical integration or expansion methods. Molenberghs and Verbeke (2005) provide an overview. We also refer to Lee, Nelder, and Pawitan (2006) for details on the computation using so-called  $h$ -likelihood methods. Most inferential approaches are based on maximum likelihood, Bayesian methods, or a variation there upon.

### 3.4 Correlation Between Two Observations Using the GLMM Framework

We will now derive a general formula for the correlation between two observations, within the GLMM framework. In the spirit of (1), and with notation consistent with Section 3.3, we can write the general model as:

$$Y_{pdt} = \mu_{pdt}^{PDT} + \varepsilon_{pdt}, \quad (14)$$

where

$$\mu_{pdt}^{PDT} = \mu_{pdt}^{PDT}(\eta_{pdt}) = h(\mathbf{x}'_{pdt}\boldsymbol{\beta} + \mathbf{z}'_{pdt}\mathbf{b}_{pdt}). \quad (15)$$

We group the errors  $\varepsilon_{pdt}$  into a vector  $\boldsymbol{\varepsilon}_p$ , with variance-covariance matrix  $\Sigma_p$ . Further, consistent with earlier notation,  $\mathbf{x}'_{pdt}$  and  $\mathbf{z}'_{pdt}$  are vectors of fixed-effects and random-effects covariates, respectively;  $\boldsymbol{\beta}$  is a vector of fixed-effects parameters and  $\mathbf{b}_{pdt}$  is a vector of random effects, assumed to be zero-mean normally distributed with variance-covariance matrix  $H$ .

It is useful in what follows to decompose  $\Sigma_p$  as:

$$\Sigma_p = \Phi^{\frac{1}{2}} A_p^{\frac{1}{2}} R_p A_p^{\frac{1}{2}} \Phi^{\frac{1}{2}}, \quad (16)$$

where  $\Phi$  is a diagonal matrix with the overdispersion parameters along the diagonal. In case there are no overdispersion parameters,  $\Phi$  is set equal to the identity matrix. Further,  $R_p$  is the correlation matrix, and  $A_p$  is a diagonal matrix containing the variances following from the generalized linear model specification of  $Y_{pdt}$  given the random effects  $\mathbf{b}_{pdt} = \mathbf{0}$ , i.e., with diagonal elements  $v(\mu_{pdt}^{PDT} | \mathbf{b}_{pdt} = \mathbf{0})$ .

Model (15) allows for a variety of distributions for the outcome variable and a wide range of link functions, while the modeler has the freedom to include or leave out serial correlation. To calculate correlation  $\text{Corr}(Y_{pdt}, Y_{p'd't'})$ , we first derive a general expression for the variance:

$$\text{Var}(Y_{pdt}) = \text{Var}(\mu_{pdt}^{PDT} + \varepsilon_{pdt}) = \text{Var}(\mu_{pdt}^{PDT}) + \text{Var}(\varepsilon_{pdt}) + 2\text{Cov}(\mu_{pdt}^{PDT}, \varepsilon_{pdt}). \quad (17)$$

It can be shown that (Molenberghs and Verbeke 2005)

$$\text{Cov}(\mu_{pdt}^{PDT}, \varepsilon_{pdt}) = \text{Cov}[E(\mu_{pdt}^{PDT} | \mathbf{b}_{pdt}), E(\varepsilon_{pdt} | \mathbf{b}_{pdt})] + E(\text{Cov}(\mu_{pdt}^{PDT}, \varepsilon_{pdt} | \mathbf{b}_{pdt})) = 0,$$

since the first term is zero and the second term equals  $E[E(\mu_{pdt}^{PDT} - E(\mu_{pdt}^{PDT}))(\varepsilon_{pdt})|\mathbf{b}_{pdt}] = 0$  as  $\mu_{pdt}^{PDT}$  is constant when conditioning on  $\mathbf{b}_{pdt}$ . For the first term in (17), we have:

$$\begin{aligned}
\text{Var}(\mu_{pdt}^{PDT}) &= \text{Var}[\mu_{pdt}^{PDT}(\eta_{pdt})] = \text{Var}[\mu_{pdt}^{PDT}(\mathbf{x}'_{pdt}\boldsymbol{\beta} + \mathbf{z}'_{pdt}\mathbf{b}_{pdt})] \\
&\cong \left( \frac{\partial \mu_{pdt}^{PDT}}{\partial \mathbf{b}_{pdt}} \bigg|_{\mathbf{b}_{pdt}=0} \right) \text{Var}(\mathbf{b}_{pdt}) \left( \frac{\partial \mu_{pdt}^{PDT}}{\partial \mathbf{b}_{pdt}} \bigg|_{\mathbf{b}_{pdt}=0} \right)' \\
&\cong \left( \frac{\partial \mu_{pdt}^{PDT}}{\partial \eta_{pdt}} \frac{\partial \eta_{pdt}}{\partial \mathbf{b}_{pdt}} \bigg|_{\mathbf{b}_{pdt}=0} \right) H \left( \frac{\partial \mu_{pdt}^{PDT}}{\partial \eta_{pdt}} \frac{\partial \eta_{pdt}}{\partial \mathbf{b}_{pdt}} \bigg|_{\mathbf{b}_{pdt}=0} \right)' \\
&\cong \Delta_{pdt} \mathbf{z}'_{pdt} H \mathbf{z}_{pdt} \Delta'_{pdt}, \tag{18}
\end{aligned}$$

where  $\Delta_{pdt} = \frac{\partial \mu_{pdt}^{PDT}}{\partial \eta_{pdt}} \bigg|_{\mathbf{b}_{pdt}=0}$ . Note that the above derivation is based on the delta method (Welsh 1996).

For the second term in (17), we have:

$$\text{Var}(\varepsilon_{pdt}) = \text{Var}[E(\varepsilon_{pdt}|\mathbf{b}_{pdt})] + E[\text{Var}(\varepsilon_{pdt}|\mathbf{b}_{pdt})] = E[\text{Var}(\varepsilon_{pdt}|\mathbf{b}_{pdt})] = \left( \Phi^{\frac{1}{2}} \Sigma \Phi^{\frac{1}{2}} \right)_{pdt}, \tag{19}$$

If the canonical link is used, we have  $A_p = \Delta_p$  and then (17) becomes

$$\text{Var}(\mathbf{Y}_p) \cong \Delta_p Z_p D Z'_p \Delta'_p + \Phi^{\frac{1}{2}} \Delta_p^{\frac{1}{2}} R_p \Delta_p^{\frac{1}{2}} \Phi^{\frac{1}{2}}. \tag{20}$$

To determine  $\text{Corr}(Y_{pdt}, Y_{p'd't'})$ , we still need to calculate  $\text{Cov}(Y_{pdt}, Y_{p'd't'})$ . Similar to the above, we have that  $\text{Cov}(\mu_{pdt}^{PDT}, \varepsilon_{p'd't'}) = \text{Cov}(\varepsilon_{pdt}, \mu_{p'd't'}^{PDT}) = 0$ . Therefore, we only need to derive  $\text{Cov}(\mu_{pdt}^{PDT}, \mu_{p'd't'}^{PDT})$ :

$$\begin{aligned}
\text{Cov}(Y_{pdt}, Y_{p'd't'}) &= \text{Cov}(\mu_{pdt}^{PDT}, \mu_{p'd't'}^{PDT}) \\
&= \text{Cov}[\mu_{pdt}^{PDT}(\mathbf{x}'_{pdt}\boldsymbol{\beta} + \mathbf{z}'_{pdt}\mathbf{b}_{pdt}), \mu_{p'd't'}^{PDT}(\mathbf{x}'_{p'd't'}\boldsymbol{\beta} + \mathbf{z}'_{p'd't'}\mathbf{b}_{p'd't'})] \\
&\cong \left( \frac{\partial \mu_{pdt}^{PDT}}{\partial \mathbf{b}_{pdt}} \bigg|_{\mathbf{b}_{pdt}=0} \right) \text{Cov}(\mathbf{b}_{pdt}, \mathbf{b}_{p'd't'}) \left( \frac{\partial \mu_{p'd't'}^{PDT}}{\partial \mathbf{b}_{p'd't'}} \bigg|_{\mathbf{b}_{p'd't'}=0} \right)' \\
&\cong \left( \frac{\partial \mu_{pdt}^{PDT}}{\partial \eta_{pdt}} \frac{\partial \eta_{pdt}}{\partial \mathbf{b}_{pdt}} \bigg|_{\mathbf{b}_{pdt}=0} \right) \text{Cov}(\mathbf{b}_{pdt}, \mathbf{b}_{p'd't'}) \left( \frac{\partial \mu_{p'd't'}^{PDT}}{\partial \eta_{p'd't'}} \frac{\partial \eta_{p'd't'}}{\partial \mathbf{b}_{p'd't'}} \bigg|_{\mathbf{b}_{p'd't'}=0} \right)' \\
&\cong \Delta_{pdt} \mathbf{z}'_{pdt} \text{Cov}(\mathbf{b}_{pdt}, \mathbf{b}_{p'd't'}) \mathbf{z}_{p'd't'} \Delta'_{p'd't'}. \tag{21}
\end{aligned}$$

The covariances  $\text{Cov}(b_{pdt}, b_{p'd't'})$  depend on which of the random effects are common when correlating  $Y_{pdt}$  and  $Y_{p'd't'}$ . Using (20) and (21), we can calculate the correlation for any given situation, for any give GLMM. In the next section, we will derive the correlation for the case of binary data with random effects and without serial correlation.

Note that, in the special case of Gaussian outcomes, (20) simply reduces to  $\text{Var}(\mathbf{Y}_p) = \mathbf{Z}_p \mathbf{H} \mathbf{Z}_p' + \mathbf{R}_p$ .

The above calculations are general, in the sense that the variances in (20) and covariances in (21) allow for the flexibe calculation of correlation coefficients, with (1) certain facets the same or different; (2) certain facets fixed (conditioned upon) or random; (3) correction for the presence of such fixed effects as treatment, time, country, baseline value, etc.; (4) for normally distributed outcomes based on linear mixed models or for binary data, count data, and other non-Gaussian data, using generalized linear mixed models. The price to pay is twofold. First, expressions (20) and (21) are approximate, except in the normal case and (2), related to the previous point, these expressions do not have the intuitive variance-component structure, or even ‘averaging’ structure, of classical reliability and generalizability coefficients. However, all classical expressions follow as special cases. In this sense, our framework allows for the calculation of conventional reliability and generalizability coefficients, their extensions to the non-normal case based on data from clinical trials or other data with measurements that are *a priori* not parallel, and even correlation coefficients that do not have a generalizability intepretation, but may be useful for other purposes.

## 4 Data Analysis

Let us now apply the concepts of reliability and generalizability to the pooled data described in Section 2. We will investigate the impact of ‘country’ on measurement error. Note that country can be seen as either a facet or an object of measurement. The generality of our approach allows for both views. Evidently, ‘country’ is of interest for this particular study, but the reader can easily substitute it with other variables, subject to his/her study of interest.



To illustrate the methods and underscore generality, we will consider country in five different roles, the analysis of which is all within reach by way of the modeling ideas developed in this manuscript. First, we will assess the overall reliability for a dichotomized version of CGI response, ignoring country effects. Second, country effects will be extracted by including country as a fixed effect into the model. Third, we will investigate the impact of country on reliability through application of the same model to each country separately. Fourth, we will also study the impact of a single country on overall reliability by leave-one-out ideas, i.e., by omitting one country at a time. Fifth, we will assess the overall impact of country via generalizability theory. Note that ‘country’ did not feature explicitly in Section 3. However, the methodology is general and the facets generic. They can be replaced with those relevant in a particular case study.

#### 4.1 Overall Reliability of CGI

First, we apply a simple random-intercept model, combined with fixed effects for treatment, time and their interaction. Hence, country does play a role in this analysis. With the logit link, (14) becomes:

$$Y_{pdt} = \frac{\exp(\mu + b_p + \mu_d^D + \mu_t^T + \mu_{dt}^{DT})}{1 + \exp(\mu + b_p + \mu_d^D + \mu_t^T + \mu_{dt}^{DT})} + \varepsilon_{pdt}, \quad (22)$$

where  $\mu_d^D$ ,  $\mu_t^T$ , and  $\mu_{dt}^{DT}$  denote the fixed effects for day, treatment, and their interaction, respectively, and  $b_p$  represents the random patient effect.

The overall correlation of observations within the same subject, on the same treatment, but on different time points, and conditioning on treatment and time points, can be expressed as  $\text{Corr}(Y_{pdt}, Y_{pd't} \mid t, d, d')$ . In this model, we have  $Z = \mathbf{1}$  and  $H = \sigma_p^2$ , a scalar representing the variance of the random intercept, and since (22) does not include serial correlation we have that  $R_p = I$ . It is therefore possible to show that the variance covariance matrix (20) reduces to

$$\text{Var}(Y_p) \cong \Delta_p(\sigma_p^2 J) \Delta_p' + \Phi \Delta_p = \Delta_p(\sigma_p^2 J + \Phi \Delta_p^{-1}) \Delta_p',$$

where  $J$  is a rectangular matrix of ones. Furthermore,  $\Delta_p$  is a diagonal matrix with  $V_{pdt}(0)$  as diagonal elements, where the variance function  $V_{pdt}(0) = \mu_{pdt}^{PDT} \big|_{b_{pdt}=0} (1 - \mu_{pdt}^{PDT} \big|_{b_{pdt}=0})$ , and

therefore we have

$$\text{Var}(Y_{pdt}) \cong \text{diag}(V_{pdt}(0))[\sigma_p^2 J + \Phi \text{diag}(V_{pdt}(0))^{-1}] \text{diag}(V_{pdt}(0)), \quad (23)$$

$$\text{Cov}(Y_{pdt}, Y_{pd't}) \cong \text{diag}(V_{pdt}(0))[\sigma_p^2 J] \text{diag}(V_{pd't}(0)). \quad (24)$$

Based on (23) and (24), we can determine a first-order approximation of the marginal correlation between time point  $d$  and  $d'$ , which is the intraclass correlation coefficient of reliability:

$$\rho = \text{Corr}(Y_{pdt}, Y_{pd't}) = \frac{\sigma_1^2 \sqrt{V_{pdt}(0) V_{pd't}(0)}}{\sqrt{[\Phi_{pdt} + V_{pdt}(0) \sigma_2^2] \cdot [\Phi_{pd't} + V_{pd't}(0) \sigma_2^2]}}, \quad (25)$$

where  $\sigma_1^2$  represents the covariance between the random effects and  $\sigma_2^2$  is the variance resulting from the random effects. In this model,  $\sigma_1^2 = \sigma_2^2 = \sigma_p^2$  since all other covariates are fixed effects.

The delta method can be usefully applied to estimate the standard error:

$$\begin{aligned} \frac{\partial \rho}{\partial(\boldsymbol{\beta}, \boldsymbol{\lambda})} &= \left( \frac{\partial(\boldsymbol{\eta}, \boldsymbol{\sigma}^2)}{\partial(\boldsymbol{\beta}, \boldsymbol{\lambda})} \right) \left( \frac{\partial(V_{pdt}(0), V_{pd't}(0), \sigma_1^2, \sigma_2^2, \phi)}{\partial(\boldsymbol{\eta}, \boldsymbol{\sigma}^2)} \right), \\ &\times \left( \frac{\partial \rho}{\partial(V_{pdt}(0), V_{pd't}(0), \sigma_1^2, \sigma_2^2, \phi)} \right). \end{aligned} \quad (26)$$

Explicit expressions for the various components follow from straightforward linear algebra, as sketched in Appendix A. The SAS V9.1 procedure GLIMMIX was used to estimate  $\Phi$ ,  $\sigma_p^2$ , and  $V_{pdt}$ . Details on the SAS implementations are provided in Appendix B. The reader interested in more ample details on the SAS implementations and output, can obtain such from the authors, upon simple request. Table 1(a) summarizes the results.

In case of continuous data, a single-measure overall intraclass correlation coefficient reliability would have been obtained (Vangeneugden *et al* 2005). Here, for the binary data case, a separate intraclass coefficient of reliability is produced for each treatment group and each time point. From Table 1(a), we observe that the correlation is somewhat higher in the risperidone arm and that the correlation between week 1 and other time points is lower than the correlation between any two other time points that do not involve week 1. This non-constancy is, of course, not particular to this example but results from the non-Gaussian nature of the outcome.

## 4.2 Overall Reliability of CGI Response Adjusting for Country

In Section 4.1, only treatment, time, and their interaction were included. Now, we will include countries as fixed effects, which will result in intraclass coefficients of reliability per treatment, time point, and country combination. Hence, country-specific analyses result. We will not present all coefficients but merely present the coefficients for one country, the U.S.A., in Table 1(b). Additionally, we list the ICC of reliability between weeks 6 and 8 in the risperidone group for all countries in Table 2. The results for the U.S.A. are consistent with the overall results, and when we investigate the correlation between weeks 6 and 8 in the risperidone group, we observe from column 3 in Table 2 that the ICC is rather stable across countries, the lowest correlation being for Austria (0.65, s.e. 0.09) and the highest for the U.S.A., Sweden, and Spain (0.78, s.e. 0.02).

## 4.3 Overall Reliability of CGI by Country and Impact on Overall Reliability by Leaving Out a Country

When we apply the model to each country separately, we observe that the model did not always converge and estimates were less stable, especially and not surprisingly, in countries with few patients. Patients included in Finland had data up to week 6 only (Hoyberg *et al* 1993). The results are summarized in the third column of Table 2. A different way of investigating impact of country on reliability is by leaving out one country at a time. This is slightly less conventional from a classical generalizability standpoint, but it is a useful analysis to assess how much a given country can weigh in on the analyses. If the overall reliability increases, this would provide evidence for a poor reliability in this specific country. The results are summarized in the fifth column of Table 2. Note that the impact was low for all countries, again suggesting that reliability is relatively consistent across countries.

#### 4.4 Estimating Impact of Country From Generalizability Theory

Subgroup analysis by country as shown in the previous two sections can be enlightening. Now, we want to quantify their effect on measurement error and calculate a generalizability coefficient, thereby generalizing results across countries. We will add a random effect for country to the previous model, so that we have a model with time, treatment, and their interaction as fixed effects, and further country, indexed by  $c$ , and patient as random effects:

$$Y_{pdtc} = \frac{\exp(\mu + b_p + \mu_d^D + \mu_t^T + \mu_{dt}^{DT} + b_c)}{1 + \exp(\mu + b_p + \mu_d^D + \mu_t^T + \mu_{dt}^{DT} + b_c)} + \varepsilon_{pdtc}. \quad (27)$$

From (27) we can calculate the overall test-retest reliability coefficient as in Section 4.1, but this time accounting for country as a random effect instead of extracting it as a fixed effect. Then,  $\sigma_1^2 = \sigma_2^2 = \sigma_p^2 + \sigma_c^2$  in (25). Table 1(c) shows that the results are consistent with the overall reliability coefficients.

This test-retest reliability coefficient for any given country and time point follows directly from analyzing the clinical trial, similar to generalizability coefficients that are computed after design and analysis of a G-study. In the spirit of D-studies, we can also generalize across countries. Indeed, although patients are nested within country in a clinical-trial setting, we assume, by way of a thought experiment, that patients could switch from one country to another, with the aim to evaluate the impact of country. We then have that  $\sigma_1^2 = \sigma_p^2$  and  $\sigma_2^2 = \sigma_p^2 + \sigma_c^2$ , needed to calculate  $\text{Corr}(Y_{pdtc}, Y_{pd'tc'})$  as in (25). Table 1(d) provides the ensuing ICC coefficients.

Thus, generalizing across time points and countries, or taking account of impact of variance of country, reduces the overall test-retest reliability approximately by 5%: for risperidone the decrease in reliability amounted to between 4–7% and for active control this was between 3–6%. In this situation, the price for setting up an international trial instead of a single country is rather small. This insight is relevant and underscores the usefulness of the thought experiment. While, again, the ‘country’ aspect will be less relevant, or even irrelevant, to the reader’s own study, our results indicate that it is possible to study the impact of generalizing over a given variable.

Evidently, the methodology can easily be extended to more complex situations including, for

example, serial correlation or random time effects but also additional variables, such as, for example, age and sex of the patient.

#### **4.5 Estimating Impact of Baseline PANSS Negative Subtotal on Reliability of CGI Response**

In the computations reported above, a relatively high generalizability coefficient suggested that country does not have an important impact on the test-retest reliability and on measurement error. We now investigate the impact of baseline PANSS Negative subtotal on measurement error. We included a random intercept for baseline PANSS Negative subtotal instead of country in Model (27). Subsequently, we derived the variance components and calculated the generalizability coefficient for baseline PANSS Negative subtotal, similar to how it was done for country. In this analysis, the reduction in generalizability coefficient was more substantial: in the risperidone group between week 6 and 8, we have that the ICC reduces from 0.55 (s.e. 0.13) to 0.39 (s.e. 0.13) when generalizing across baseline negative subtotal. Full details are given in Table 1(e). This indicates that baseline PANSS Negative subtotal reduces the test-retest reliability. A clinical explanation for this phenomenon could be that patients with a higher deficit in negative symptoms at baseline, such as poverty of speech, apathy, or emotional withdrawal, are more difficult to evaluate, resulting in higher measurement error and lower test-retest reliability. A practical conclusion would be that additional training is needed for professionals having to rate patients with a high baseline negative subtotal or, even more invasive, in the recommendation to use a different scale in this type of patients. Such conclusions usefully illustrate how the methodology can be used, not only to assess the qualitative level of generalizability, but also how such results can impact the design of future studies.

## 5 Concluding Remarks

In this paper, we have extended classical reliability measures and associated estimation procedures in four important ways. First, fully longitudinal data can be used, rather than paired measurements. Second, clinical trial data can be employed or, more generally data from other studies not expressly designed for the investigation of reliability, through adopting a modeling framework, obviating the need for parallel measurements. Third, the broad generalizability theory framework is invoked, encompassing the various classical reliability versions, such as inter-rater and test-retest reliability, and allowing for the study of such important factors' impact as day of measurement, rater, country, investigator, etc. Fourth, all calculations are conducted within the generalized linear mixed model paradigm, allowing one not only to accommodate all aforementioned aspects, but also to deal with Gaussian and non-Gaussian data alike. Specific emphasis was put on binary outcomes, but analogous computations for nominal, ordinal, or count data can be done as well. Unlike in the Gaussian case, the reliability and generalizability coefficients depend on the days, raters, countries, or whatever levels studied. This is due to the mean-variance link and the nonlinear nature of the model.

Of course, our calculations are based on a first-order approximation, the accuracy of which could be a cause of concern. Vangeneugden *et al* (2007) have studied this issue and their results are surprisingly encouraging.

We would like to emphasize that we have focused on generalizability, with reliability as a special case. This implies that we have been less concerned with agreement. While the latter concept is also very important, it falls outside the scope of the current work.

This work was motivated by and applied to data from multi-country trial data collected in patients with chronic schizophrenia. Using the generalizability framework, we were able to establish that the reliability measures are rather stable across countries, and no single country has an undue effect on the overall reliability. Country-specific reliabilities varied in a usefully narrow range.

An important conclusion, never reached before, is that the price to pay for a multi-country study,

rather than a single-country one, is a mere 5% in test-retest reliability. The ability to conduct multi-country studies is important in view of the availability of a larger pool of available patients, thereby reducing the length of the accrual period and/or increasing the sample size, and hence power.

## Acknowledgments

The authors are thankful to J&J PRD for kind permission to use their data. Financial support from the IAP research network #P6/03 of the Belgian Government (Belgian Science Policy) is gratefully acknowledged. We thank the Editor, the Associate Editor, and anonymous referees for helpful and constructive comments on earlier versions of this work.

## References

- Bell, M., Milstein, R., Beam-Goulet, J., Lysaker, P., and Cicchetti, D. (1992). The Positive and Negative Syndrome Scale and the Brief Psychiatric Rating Scale: Reliability, comparability, and predictive validity. *Journal of Nervous and Mental Disease* **180**, 723–728.
- Brennan, R.L. (1992). *Elements of Generalizability Theory*. Iowa City, IA: ACT Publications.
- Breslow, N.E. and Clayton, D.G. (1993). Approximate inference in generalized linear mixed models. *Journal of American Statistical Association* **88**, 9–25.
- Cronbach, L.J., Rajaratnam, N., and Gleser, G.C. (1963). Theory of generalizability: a liberalization of reliability theory. *British Journal of Statistical Psychology* **16**, 137–163.
- Cronbach, L.J., Gleser, G.C., Nanda, H., and Rajaratnam, N. (1972). *The Dependability of Behavioral Measurements: Theory of Generalizability for Scores and Profiles*. New York: John Wiley.
- Fleiss J.L. (1981). *Statistical Methods for Rates and Proportions*. New York: John Wiley.
- Fleiss, J.L. (1986). *Design and Analysis of Clinical Experiments*. New York: John Wiley.
- Hoyberg, O.J., Fensbo, C., Remvig, J., Lingjaerde, O., Sloth-Nielsen, M., and Salvesen, I. (1993).

- Risperidone versus perphenazine in the treatment of chronic schizophrenic patients with acute exacerbations. *Acta Psychiatrica Scandinavica* **88**, 395–402.
- Huttunen, M.O., Piepponen, T., Rantanen, H., Larmo, L., Nyholm, R., and Raitasuo, V. (1995). Risperidone versus zuclopenthixol in the treatment of acute schizophrenic episodes: a double-blind parallel-group trial. *Acta Psychiatrica Scandinavica* **91**, 271–277.
- Kay, S.R., Fiszbein, A., and Opler, L.A. (1987). The Positive and Negative Syndrome Scale (PANSS) for Schizophrenia. *Schizophrenia Bulletin* **13**, 261–276.
- Kay, S.R., Opler, L.A., and Lindenmayer, J.P. (1988). Reliability and validity of the Positive and Negative Syndrome Scale for Schizophrenics. *Psychiatric Research* **23**, 99–110.
- Lee, Y., Nelder, J.A., and Pawitan, Y. (2006). *Generalized Linear Models with Random Effects*. Boca Raton: Chapman & Hall/CRC.
- Liang, K.Y. and Zeger, S.L. (1986). Longitudinal Data Analysis Using Generalized Linear Models. *Biometrika* **73**, 13–22.
- Marder, S.R. and Meibach, R.C. (1994). Risperidone in the treatment of schizophrenia. *American Journal of Psychiatry* **151**, 825–835.
- McCullagh, P. and Nelder, J.A. (1989). *Generalized Linear Models*. London: Chapman & Hall.
- Molenberghs, G. and Verbeke, G. (2005). *Models for Discrete Longitudinal Data*. New York: Springer.
- Peralta, V. and Cuesta, M.J. (1994). Psychometric properties of the Positive and Negative Syndrome Scale (PANSS) in Schizophrenia. *Psychiatric Research* **53**, 31–40.
- Peuskens, J. and the Risperidone Study Group (1995). Risperidone in the treatment of chronic schizophrenic patients: a multinational, multicentre, double-blind, parallel-group study versus haloperidol. *British Journal of Psychiatry* **166**, 712–726.
- Shavelson, R.J., Webb, N.M., and Rowley, G.L. (1989). Generalizability theory. *American Psychologist* **44**, 922–932.
- Shrout, P.E. and Fleiss, J.L. (1979). Intraclass correlations: uses in assessing interrater reliability. *Psychological Bulletin* **86**, 420–428.



- Stratford P. (1989). Consistency or differentiating among subjects? *Physical Therapy* **69**, 299–300.
- Streiner D.L and Norman G.R. (1995). *Health Measurement Scales*. Oxford University Press.
- Vangeneugden, T., Laenen, A., Geys, H., Renard, D. and Molenberghs, G. (2004). Applying linear mixed models to estimate reliability in clinical trial data with repeated measurements. *Controlled Clinical Trials* **25**, 13–30.
- Vangeneugden, T., Laenen, A., Geys, H., Renard, D. and Molenberghs, G. (2005). Applying concepts of generalizability theory on clinical trial data to investigate sources of variation and their impact on reliability. *Biometrics* **61**, 295–304.
- Vangeneugden, T., Molenberghs, G., Laenen, A., Geys, H., Beunckens, C., and Sotto, C. (2007) Marginal correlation in longitudinal binary data based on generalized linear mixed models. *Submitted for publication*.
- Welsh, A.H. (1996) *Aspects of Statistical Inference*. New York: Wiley.

## A Explicit Expressions for Components of (26)

For notational simplicity, write  $\pi_0 \equiv V_{pd^t}(0)$  and  $\pi'_0 \equiv V_{pd'^t}(0)$ . Further, note that  $\boldsymbol{\eta} = (\eta_1 = \mathbf{x}'_1\boldsymbol{\beta}, \eta_2 = \mathbf{x}'_2\boldsymbol{\beta})'$ ,  $\boldsymbol{\sigma}^2 = (\sigma_1^2, \sigma_2^2, \phi)'$ , and  $\boldsymbol{\lambda} = (\lambda_1, \dots, \lambda_K, \phi)'$ , where the  $\lambda_k$  are the parameters figuring in the models for  $\sigma_1^2$  and  $\sigma_2^2$ . We then obtain, for the first factor on the right hand side:

$$\frac{\partial(\boldsymbol{\eta}, \boldsymbol{\sigma}^2)}{\partial(\boldsymbol{\beta}, \boldsymbol{\lambda})} = \begin{pmatrix} \mathbf{x}'_1 & 0 & 0 & 0 \\ 0 & \mathbf{x}'_2 & 0 & 0 \\ 0 & \frac{\partial\sigma_1^2}{\partial\lambda_1} & \frac{\partial\sigma_2^2}{\partial\lambda_1} & 0 \\ \vdots & \vdots & \vdots & \vdots \\ 0 & \frac{\partial\sigma_1^2}{\partial\lambda_K} & \frac{\partial\sigma_2^2}{\partial\lambda_K} & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix}.$$

The second factor equals:

$$\frac{\partial(\pi_0, \pi'_0, \sigma_1^2, \sigma_2^2, \phi)}{\partial(\eta_1, \eta_2, \sigma_1^2, \sigma_2^2, \phi)} = \left( \begin{array}{cc|c} \pi_0(1 - \pi_0) & 0 & \mathbf{0}_{2,3} \\ 0 & \pi'_0(1 - \pi'_0) & \\ \hline \mathbf{0}_{3,2} & & I_{3,3} \end{array} \right).$$

To establish the elements of the third factor, first write  $\rho = \sigma_1^2 \kappa_0^{1/2} \kappa_1^{-1/2} \kappa_2^{-1/2}$ , with  $\kappa_0 = \pi_0(1 - \pi_0)\pi'_0(1 - \pi'_0)$ ,  $\kappa_1 = \pi_0(1 - \pi_0)\sigma_2^2 + \phi$ , and  $\kappa_2 = \pi'_0(1 - \pi'_0)\sigma_2^2 + \phi$ . It then follows that

$$\begin{aligned} \frac{\partial\rho}{\partial\pi_0} &= \frac{1}{2}\sigma_1^2\kappa_0^{-1/2}\kappa_1^{-3/2}\kappa_2^{-1/2}(1 - 2\pi_0)\pi'_0(1 - \pi'_0)\phi, \\ \frac{\partial\rho}{\partial\pi'_0} &= \frac{1}{2}\sigma_1^2\kappa_0^{-1/2}\kappa_1^{-1/2}\kappa_2^{-3/2}(1 - 2\pi'_0)\pi_0(1 - \pi_0)\phi, \\ \frac{\partial\rho}{\partial\sigma_1^2} &= \kappa_0^{1/2}\kappa_1^{-1/2}\kappa_2^{-1/2}, \\ \frac{\partial\rho}{\partial\sigma_2^2} &= -\frac{1}{2}\sigma_1^2\kappa_0^{1/2}\kappa_1^{-3/2}\kappa_2^{-3/2}\{2\kappa_0\sigma_2^2 + [\pi_0(1 - \pi_0) + \pi'_0(1 - \pi'_0)]\phi\}, \\ \frac{\partial\rho}{\partial\phi} &= -\frac{1}{2}\sigma_1^2\kappa_0^{1/2}\kappa_1^{-3/2}\kappa_2^{-3/2}\{\pi_0(1 - \pi_0) + \pi'_0(1 - \pi'_0)]\sigma_2^2 + 2\phi\}. \end{aligned}$$

## B SAS Implementation

All data analyses have been conducted using the SAS procedure GLIMMIX to obtain parameter estimates and measures of precision. The correlation quantities derived have been obtained using user-defined code written in the SAS procedure IML. Here, we will provide some example code and a brief discussion. A full set of programs and output can be obtained from the authors' web pages.

First, a SAS program, using the procedure GLIMMIX, for the model of Section 4.1, with a random intercept and a scale parameter  $\Phi$ , taking into account treatment, time, as well as their interaction, is as follows:

```
ods output ParameterEstimates=datagen.model1estimates
           covparms=datagen.model1covparms
           covb=datagen.model1covb
           asycov=datagen.model1asycov;

proc glimmix data=cgi3 noclprint=26 asycov;
  class treat id  xtime_c;
  title 'Overall model with random patient to estimate
        overall reliability with scale parameter';
  model cgi_resp(event='1') = xtime_c|treat
    / covb dist=binary link=logit solution;
  random intercept / subject=id;
  random _residual_;
run;
```

The coding is self-explanatory, in the sense that the fixed-effects structure involves time, treatment, and their interaction, and a random intercept is then added. The RANDOM \_residual\_ statement ensures the scale, or overdispersion, parameter is included.

Second, to estimate the correlation and its associated standard error, SAS IML can be used. Let us exemplify this for the correlation between the first and second occasions, within the risperidone group.

```
proc iml;
  use datagen.model1estimates;
  read all var {estimate} into beta;
  use datagen.model1covparms;
  read all var {estimate} into sigma;
  use datagen.model1asycov;
  read all var {CovP1,CovP2} into asycov;
  close datagen.model1asycov;
  use datagen.model1covb;
  read all var {Col1,Col2,Col3,Col4,Col5,Col6,Col7,Col8,Col9,Col10,Col11,
               Col12,Col13,Col14,Col15,Col16,Col17,Col18} into covb;
  close datagen.model1covb;
  varint=sigma[1];
  scale=sigma[2];
  covarmatrix=block(covb,asycov);
  zero=J(18,1,0);

  *           Time      RX   Contr*Time Ris*time;
  *           I  1 2 4 6 8  C R  1 2 4 6 8  1 2 4 6 8;
  bw1ris=T({1  1 0 0 0 0  0 1  0 0 0 0 0  1 0 0 0 0});
  bw2ris=T({1  0 1 0 0 0  0 1  0 0 0 0 0  0 1 0 0 0});
  bw4ris=T({1  0 0 1 0 0  0 1  0 0 0 0 0  0 0 1 0 0});
  bw6ris=T({1  0 0 0 1 0  0 1  0 0 0 0 0  0 0 0 1 0});
  bw8ris=T({1  0 0 0 0 1  0 1  0 0 0 0 0  0 0 0 0 1});

  bw1con=T({1  1 0 0 0 0  1 0  1 0 0 0 0  0 0 0 0 0});
```

```

bw2con=T({1  0 1 0 0 0  1 0  0 1 0 0 0  0 0 0 0 0});
bw4con=T({1  0 0 1 0 0  1 0  0 0 1 0 0  0 0 0 0 0});
bw6con=T({1  0 0 0 1 0  1 0  0 0 0 1 0  0 0 0 0 0});
bw8con=T({1  0 0 0 0 1  1 0  0 0 0 0 1  0 0 0 0 0});

pilris=exp(T(bw1ris)*beta)/(1+exp(T(bw1ris)*beta));
pi2ris=exp(T(bw2ris)*beta)/(1+exp(T(bw2ris)*beta));

k0_1_2_ris=(pilris*(1-pilris)*pi2ris*(1-pi2ris));
k1_1_2_ris=pilris*(1-pilris)*varint+scale;
k2_1_2_ris=pi2ris*(1-pi2ris)*varint+scale;

* correlation following equation (23)
r1_2_ris=(sqrt(k0_1_2_ris)*varint)/(sqrt(k1_1_2_ris)*sqrt(k2_1_2_ris));

D_rho_sigma_1_1_2_ris=sqrt(k0_1_2_ris)/(sqrt(k1_1_2_ris)*sqrt(k2_1_2_ris));
D_rho_sigma_2_1_2_ris=-0.5*varint*sqrt(k0_1_2_ris)*(k1_1_2_ris*k2_1_2_ris)**(-1.5)
    *(2*k0_1_2_ris*varint
    +scale*(pilris*(1-pilris)+pi2ris*(1-pi2ris)));
D_rho_poa_1_2_ris=0.5*varint*(k0_1_2_ris*k2_1_2_ris)**(-0.5)
    *(k1_1_2_ris)**(-1.5)*(scale*pi2ris*(1-pi2ris)*(1-2*pilris));
D_rho_poa_1_2_ris=0.5*varint*(k0_1_2_ris*k1_1_2_ris)**(-0.5)
    *(k2_1_2_ris)**(-1.5)*(scale*pilris*(1-pilris)*(1-2*pi2ris));
D_rho_phi_1_2_ris=-0.5*varint*sqrt(k0_1_2_ris)
    *(k1_1_2_ris*k2_1_2_ris)**(-1.5)*(k1_1_2_ris+k2_1_2_ris);
F_ris_1_2=D_rho_poa_1_2_ris*pilris*(1-pilris)*T(bw1ris)+D_rho_poa_1_2_ris*pi2ris
    *(1-pi2ris)*T(bw2ris) ||
D_rho_sigma_1_1_2_ris+D_rho_sigma_2_1_2_ris||D_rho_phi_1_2_ris ;

```

```

se_ris_1_2=sqrt(F_ris_1_2*covarmatrix*T(F_ris_1_2));

print "Estimated Correlation matix-Risperdal" cor_ris[format=8.2];
print "Standard error correlations-Risperdal" se_ris[format=8.2];
print "Estimated Correlation matix Control" cor_con[format=8.2];
print "Standard error correlations-Control" se_con[format=8.2];
quit;

```

The IML code is a little tedious, but otherwise reasonably straightforward. Different models require slightly modified coding of the GLIMMIX procedure, while the IML code needs adaptation as well.

The model of Section 4.2 requires replacement of two statement in the GLIMMIX code:

```

class treat id xtime_c country;
model cgi_resp(event='1') = xtime_c|treat country
      / covb dist=bin link=logit solution;

```

The IML code is more extensive, since a specific contribution for each country is calculated.

The analyses of Section 4.3, by country on the one hand and using leave-on-country-out on the other hand, are done by applying macros:

```

%bycountry(inds=cgi3,land='ARG');
%countryout(land="ARG");

```

The program for Section 4.4, with country as random effect, is the same as the program for Section 4.1, i.e., the first one presented in this appendix, with simply the following statement added:

```

random intercept / subject=country;

```

in addition to the two RANDOM statements already present.

**Table 1:** ICC matrices (standard error), accounting for treatment, time and their interaction.

Standard errors are calculated from the delta method. Five different situations are reported.

Week	risperidone				active control			
	2	4	6	8	2	4	6	8
(a) Overall								
1	.52(.04)	.55(.04)	.55(.04)	.55(.04)	.42(.04)	.47(.04)	.50(.04)	.50(.04)
2	1	.74(.02)	.74(.02)	.74(.02)	1	.61(.04)	.65(.03)	.66(.03)
4		1	.78(.02)	.78(.02)		1	.72(.03)	.73(.02)
6			1	.79(.01)			1	.78(.02)
(b) By country: U.S.A.								
1	.52(.06)	.54(.06)	.54(.05)	.54(.05)	.38(.07)	.42(.07)	.46(.06)	.46(.06)
2	1	.73(.03)	.74(.03)	.74(.02)	1	.57(.06)	.62(.05)	.63(.05)
4		1	.77(.02)	.77(.02)		1	.69(.04)	.70(.04)
6			1	.78(.02)			1	.76(.02)
(c) Country as random effect: U.S.A.								
1	.53(.05)	.55(.05)	.56(.05)	.56(.05)	.40(.06)	.45(.06)	.48(.05)	.48(.05)
2	1	.74(.03)	.75(.02)	.75(.02)	1	.59(.05)	.64(.04)	.65(.04)
4		1	.78(.02)	.78(.02)		1	.71(.03)	.72(.03)
6			1	.79(.02)			1	.77(.02)
(d) Generalized across countries: U.S.A.								
1	.49(.05)	.51(.05)	.51(.05)	.51(.04)	.37(.05)	.41(.05)	.44(.05)	.45(.05)
2	1	.68(.03)	.69(.03)	.69(.03)	1	.55(.05)	.59(.04)	.60(.04)
4		1	.72(.03)	.72(.03)		1	.65(.04)	.66(.03)
6			1	.72(.03)			1	.71(.03)
(e) Generalized across baseline negative symptoms								
1	.37(.13)	.38(.13)	.39(.13)	.39(.13)	.29(.10)	.32(.11)	.35(.12)	.35(.12)
2	1	.51(.18)	.52(.18)	.52(.18)	1	.43(.15)	.46(.16)	.46(.16)
4		1	.54(.18)	.54(.18)		1	.50(.17)	.51(.17)
6			1	.55(.19)			1	.54(.18)

**Table 2:** *Reliability by country and impact of country on overall reliability table. ICC  $\rho$  (standard error) between Week 6 and 8 in risperidone, with (1) country as fixed effect, (2) country-specific analyzes, and (3) a given country omitted. (NA: not available by lack of data.)*

Country	Number of patients	Country as fixed effect	By country	Omitting a given country
Argentina	31	0.76 (0.04)	NA	0.78 (0.02)
Austria	29	0.65 (0.09)	0.02 (0.04)	0.78 (0.01)
Belgium	26	0.76 (0.04)	NA	0.78 (0.01)
Brazil	44	0.73 (0.05)	0.54 (0.14)	0.79 (0.01)
Canada	44	0.77 (0.02)	0.76 (0.10)	0.79 (0.01)
Denmark	47	0.77 (0.02)	0.65 (0.09)	0.80 (0.01)
Spain	32	0.78 (0.02)	0.88 (0.07)	0.79 (0.01)
Finland	71	0.66 (0.07)	NA	0.79 (0.01)
France	92	0.77 (0.02)	0.40 (0.11)	0.81 (0.01)
Great Britain	21	0.77 (0.03)	0.91 (0.05)	0.78 (0.01)
Germany	25	0.73 (0.06)	NA	0.78 (0.01)
Italy	39	0.70 (0.07)	NA	0.77 (0.02)
Mexico	36	0.76 (0.03)	0.92 (0.06)	0.78 (0.02)
Netherlands	17	0.74 (0.06)	0.71 (0.37)	0.78 (0.01)
Norway	37	0.71 (0.06)	0.91 (0.04)	0.78 (0.01)
South Africa	79	0.71 (0.05)	0.80 (0.09)	0.78 (0.02)
Sweden	30	0.78 (0.02)	0.94 (0.03)	0.78 (0.01)
U.S.A.	122	0.78 (0.02)	0.75 (0.04)	0.79 (0.02)