

Contents

1	Introduction to Method Comparison Studies	9
1.1	Agreement	9
1.2	Purposes of MCS	9
1.3	Method Comparison Studies	10
1.4	Discussion on Method Comparison Studies	10
1.5	Indications on how to deal with outliers in Bland Altman plots	11
1.5.1	Agreement	11
1.6	Methods of assessing agreement	12
1.6.1	Equivalence and Interchangeability	12
1.7	Introductory Definitions	13
1.8	Agreement Criteria	16
1.9	Roy's Approach	17
1.10	LME Model Specification	19
1.10.1	Variance Covariance Matrices	20
1.11	Likelihood Ratio Tests	21
1.11.1	Statistical Assumptions for Likelihood Ratio Tests	22
1.11.2	Nesting: Model Selection Using Likelihood Ratio Tests	23
1.11.3	Relevance of Estimation Methods	23
1.11.4	Akaike Information Critierion	24
2	Linear Mixed effects Models	25

2.1	Case Deletion Diagnostics for LME models	25
2.2	Model Terms (Roy 2009)	26
2.3	BXC - Model Terms	27
2.4	LME	28
2.5	Remarks	29
3	LME Likelihood	30
3.1	PRESS	30
3.2	One Way ANOVA	30
3.2.1	Page 448	30
3.2.2	Page 448- simple example	31
3.3	Sampling	32
3.4	Conclusion	32
3.4.1	EBLUPS-Diagnostics for Random Effects	33
4	General Appendices	34
4.0.1	Extending deletion diagnostics to LMEs	34
4.1	Unknown Material	35
4.1.1	Estimation	36
4.1.2	Zewotir-Cook's Distance	36
4.1.3	Leverage	37
4.1.4	Local Influence	38
4.2	ICC, Reproducibility Index and Passing-Bablok	45
4.2.1	Intraclass Correlation Coefficient	45
4.2.2	Passing and Bablok (1983)	45
4.2.3	Lin's Reproducibility Index	46
4.3	Repeated Measurements	46
4.4	Linnet - References	46
4.5	Lewis Conversion	47
4.6	Likelihood ratio test	47

4.7	RSquared for LME models	49
4.8	Remarks on the Multivariate Normal Distribution	50
4.8.1	Lin's Reproducibility Index	51
4.9	Measurement Error Models	51
4.9.1	The Problem of Identifiability	51
4.9.2	Identifiability	52
4.10	Carstensen Model (mir model)	52
5	Bradley Blackwood	55
5.1	Bartko's Bradley-Blackwood Test	55
5.2	Bradley-Blackwood Test (Kevin Hayes Talk)	56
5.3	Simple Linear Regression	57
5.4	Constant and Proportional Bias	59
5.5	Conclusions about Existing Methodologies	59
5.6	A regression based approach based on Bland Altman Analysis	61
5.7	The MCR R pacakge - Regression Techniques for MCS	61
5.8	Implementation of Deming Regression with Rs	62
5.9	KP	62
6	Residual Diagnostics	63
6.0.1	Case-Deletion Diagnostics	64
6.1	Cook's Distance	65
6.2	Haslett Dillane Hayes	65
6.3	Demidenk Case Deletion Diagnostics	65
6.4	Cooks's Distance - Implementation with R	66
6.5	Influence measures using R	66
6.6	LME diagnostic measures	66
6.6.1	Andrews-Pregibon statistic	66
6.6.2	Cook's Distance	66
6.6.3	Variance Ratio	67

6.6.4	Cook-Weisberg statistic	67
6.6.5	Andrews-Pregibon statistic	67
6.7	Two-tailed testing	68
6.8	One Tailed Testing	68
6.9	Enabling One Tailed Testing	68
6.10	Profile Likelihood	69
6.11	Implementation of PL Confidence Intervals	69
6.12	Zewotir: Computation and Notation	70
6.13	Haslett Hayes	70
6.14	Confounded Residuals	70
7	Fitting LME Models	72
7.1	Relevance of Roy's Methodology	73
7.2	Interaction Terms in Model	73
7.3	Difference Variance further to Carstensen	73
7.4	Why use LMEs for Method Comparison?	74
7.5	Definition of Replicate measurements	75
7.5.1	Exchangeable measurements	75
7.5.2	Linked measurements	75
7.5.3	Replicate measurements in ARoy2009's paper	76
7.5.4	Random effects	77
7.6	Model for replicate measurements	77
8	BA99	79
8.1	Regression-based Limits of Agreement	79
8.2	Steps of Structural Equation modelling	80
9	Appendices 1	81
9.1	Model Terms (ARoy2009 2009)	81

10 Augmented GLMs	82
10.1 Augmented GLMs	82
10.1.1 The Augmented Model Matrix	83
10.2 Algorithms : ML v REML	83
10.3 Estimation of random effects	84
10.4 Covariance Parameters	84
10.4.1 Methods and Measures	84
10.5 Haslett's Analysis	85
10.6 Computation and Notation	85
11 Generalized linear models	86
11.1 Generalized Linear model	86
11.2 Generalized Model(GzLM)	86
11.2.1 What is a GzLM	87
11.2.2 GzLM Structure	87
11.2.3 Link Function	87
11.2.4 Canonical parameter	87
11.2.5 Dispersion parameter	87
11.2.6 Iteratively weighted least square	88
11.2.7 Residual Components	88
11.3 Generalized linear mixed models	88
11.4 Assessment of Agreements in Linear and Generalized Linear Mixed Models	89
11.5 Random Effects and MCS	90
11.5.1 Random coefficient growth curve model	90
11.6 Random effects Model	90
11.6.1 Myers Random Effects Model	91
11.6.2 Random Effects Modelling	91
11.7 Other Approaches : Marginal Modelling	92
11.8 Other Approaches	92

Bibliography	92
------------------------	----

Chapter 1

Introduction to Method Comparison Studies

1.1 Agreement

- The FDA define precision as the *closeness of agreement* (degree of scatter) between a series of measurements obtained from multiple sampling of the same homogeneous sample under prescribed conditions.
- **Barnhart** describes precision as being further subdivided as either within-run, intra-batch precision or repeatability (which assesses precision during a single analytical run), or between-run, inter-batch precision or repeatability (which measures precision over time).

1.2 Purposes of MCS

The question being answered is not always clear, but is usually expressed as an attempt to quantify the agreement between two methods (Bland and Altman 1995)

Some lack of agreement between different methods of measurement is inevitable. What matters is the amount by which they disagree. we want to know by how much the new method is likely to differ from the old, so that it is not enough to cause problems in the mathematical interpretation we can

preplace the old method by the new, or even use the two interchangeably.

It often happens that the same physical and chemical property can be measured in different ways. For example, one can determine sodium in serum by flame atomic emission spectroscopy or by isotops dilution mass spectroscopy. The question arises as to whcih methd is better (Mandel 1991)

In areas of inter-laboratory quality control, method comparisons, assay validations and individual bio-equivalence, etc, the agree between observations and target (reference) value is of interest (lin 2002)

The purpose of comparing two methods of measurement of a continuous biological variable is to uncover systematic differences, not to point to similarities. (ludbrook 1997)

In the pharmaceutical industry, measurement methods that measure the quantity of prdocuts are regulated. The FDA (U.S. Food and Drug Administration) requires that the manufacturer show equivalency prior to approving the new or alternatice method in quality control (Tan & Inglewicz ,1999)

1.3 Method Comparison Studies

Agreement between two methods of clinical measurement can be quantified using the differences between observations made using the two methods on the same subjects. The 95% limits of agreement, estimated by mean difference \pm 1.96 standard deviation of the differences, provide an interval within which 95% of differences between measurements by the two methods are expected to lie.

1.4 Discussion on Method Comparison Studies

The need to compare the results of two different measurement techniques is common in medical statistics.

In particular, in medicine, new methods or devices that are cheaper, easier to use, or less invasive, are routinely developed. Agreement between a new method and a traditional reference or gold standard must be evaluated before the new one is put into practice. Various methodologies have been proposed for this purpose in recent years.

1.5 Indications on how to deal with outliers in Bland Altman plots

We wish to determine how outliers should be treated in a Bland Altman Plot

In their 1983 paper they merely state that the plot can be used to 'spot outliers'.

In their 1986 paper, Bland and Altman give an example of an outlier. They state that it could be omitted in practice, but make no further comments on the matter.

In Bland and Altmans 1999 paper, we get the clearest indication of what Bland and Altman suggest on how to react to the presence of outliers. Their recommendation is to recalculate the limits without them, in order to test the difference with the calculation where outliers are retained.

The span has reduced from 77 to 59 mmHg, a noticeable but not particularly large reduction.

However, they do not recommend removing outliers. Furthermore, they say:

We usually find that this method of analysis is not too sensitive to one or two large outlying differences. We ask if this would be so in all cases. Given that the limits of agreement may or may not be disregarded, depending on their perceived suitability, we examine whether it would be possible that the deletion of an outlier may lead to a calculation of limits of agreement that are usable in all cases?

Should an Outlying Observation be omitted from a data set? In general, this is not considered prudent. Also, it may be required that the outliers are worthy of particular attention themselves.

Classifying outliers and recalculating We opted to examine this matter in more detail. The following points have to be considered

how to suitably identify an outlier (in a generalized sense)

Would a recalculation of the limits of agreement generally results in a compacted range between the upper and lower limits of agreement?

1.5.1 Agreement

Bland and Altman (1986) define Perfect agreement as 'The case where all of the pairs of rater data lie along the line of equality'. The Line of Equality is defined as the 45 degree line passing through the

origin, or $X=Y$ on a XY plane.

1.6 Methods of assessing agreement

1. Pearson's Correlation Coefficient
2. Intraclass correlation coefficient
3. Bland Altman Plot
4. Bartko's Ellipse (1994)
5. Blackwood Bradley Test
6. Lin's Reproducibility Index
7. Luiz Step function

Bland and Altman attend to the issue of repeated measures in 1996.

Repeated measurements on several subjects can be used to quantify measurement error, the variation between measurements of the same quantity on the same individual.

Bland and Altman discuss two metrics for measurement error; the within-subject standard deviation, and the correlation coefficient.

The above plot incorporates both the conventional limits of agreement (the inner pair of dashed lines), the 't' limits of agreement (the outer pair of dashed lines) centred around the inter-method bias (indicated by the full line). This plot is intended for expository purposes only, as the sample size is small.

1.6.1 Equivalence and Interchangeability

Limits of agreement are intended to analyse equivalence. How this is assessed is the considered judgement of the practitioner. In Bland and Altman (1986) an example of good agreement is cited. For two methods of measuring 'oxygen saturation', the limits of agreement are calculated as (-2.0,2.8). A practitioner would ostensibly find this to be sufficiently narrow.

If the limits of agreement are not clinically important, which is to say that the differences tend not to be substantial, the two methods may be used interchangeably. Dunn (2002) takes issue with the notion of ‘equivalence’, remarking that while agreement indicated equivalence, equivalence does not reflect agreement.

1.7 Introductory Definitions

To illustrate the characteristics of a typical method comparison study consider the data in Table I, taken from Grubbs (1973). In each of twelve experimental trials a single round of ammunition was fired from a 155mm gun, and its velocity was measured simultaneously (and independently) by three chronographs devices, referred to here as ‘Fotobalk’, ‘Counter’ and ‘Terma’.

Round	Fotobalk [F]	Counter [C]	Terma [T]
1	793.8	794.6	793.2
2	793.1	793.9	793.3
3	792.4	793.2	792.6
4	794.0	794.0	793.8
5	791.4	792.2	791.6
6	792.4	793.1	791.6
7	791.7	792.4	791.6
8	792.3	792.8	792.4
9	789.6	790.2	788.5
10	794.4	795.0	794.7
11	790.9	791.6	791.3
12	793.5	793.8	793.5

Table 1.7.1: Measurement of the three chronographs (Grubbs 1973)

An important aspect of the these data is that all three methods of measurement are assumed to have an attended measurement error, and the velocities reported in Table I can not be assumed to be ‘true values’ in any absolute sense. For expository purposes only the first two methods ‘Fotobalk’ and ‘Counter’ will enter in the immediate discussion.

While lack of agreement between two methods is inevitable, the question , as posed by Altman and Bland (1983), is ‘do the two methods of measurement agree sufficiently closely?’

A method of measurement should ideally be both accurate and precise. An accurate measurement methods shall give a result close to the ‘true value’. Precision of a method is indicated by how tightly clustered its measurements are around their mean measurement value.

A precise and accurate method should yield results consistently close to the true value. However a method may be accurate, but not precise. The average of its measurements is close to the true value, but those measurements are highly dispersed. Conversely an inaccurate method may be quite precise , as it consistently indicates the same level of inaccuracy.

The tendency of a method of measurement to consistently give results above or below the true value is a source of systematic bias. The lesser the systematic bias, the greater the accuracy of the method.

In the context of the agreement of two methods, there is also a tendency of one measurement method to consistently give results above or below the other method. Lack of agreement is a consequence of the existence of ‘inter-method bias’. For two methods to be considered in good agreement, the inter-method bias should be in the region of zero.

A simple estimation of the inter-method bias can be calculated using the differences of the paired measurements. The data in Table 1.2 are a good example of possible inter-method bias; the ‘Fotobalk’ consistently recording smaller velocities than the ‘Counter’ method. Consequently there is lack of agreement between the two methods.

Round	Fotobalk (F)	Counter (C)	F-C
1	793.80	794.60	-0.80
2	793.10	793.90	-0.80
3	792.40	793.20	-0.80
4	794.00	794.00	0.00
5	791.40	792.20	-0.80
6	792.40	793.10	-0.70
7	791.70	792.40	-0.70
8	792.30	792.80	-0.50
9	789.60	790.20	-0.60
10	794.40	795.00	-0.60
11	790.90	791.60	-0.70
12	793.50	793.80	-0.30

Table 1.7.2: Difference between Fotobalk and Counter measurements

The absence of inter-method bias by itself is not sufficient to establish whether two measurement methods agree or not. These methods must also have equivalent levels of precision. Should one method yield results considerably more variable than that of the other, they can not be considered to be in agreement.

Therefore a methodology must be introduced that allows an analyst to estimate the inter-method bias, and to compare the precision of both methods of measurement.

1.7.1 Agreement Criteria

Roy's method considers two methods to be in agreement if three: no significant bias, i.e. the difference between the two mean readings is not "statistically significant", high overall correlation coefficient, the agreement between the two methods by testing their repeatability coefficients. Two methods of measurement can be said to be in agreement if there is no significant difference between in three key respects.

Roy additionally uses the overall correlation coefficient to provide extra information about the comparison, with a minimum of 0.82 being required. Firstly, there is no inter-method bias between the two methods, i.e. there is no persistent tendency for one method to give higher values than the other. Secondly, both methods of measurement have the same within-subject variability. In such a case the variance of the replicate measurements would be consistent for both methods. Lastly, the methods have equal between-subject variability. Put simply, for the mean measurements for each case, the variances of the mean measurements from both methods are equal. Lack of agreement can arise if there is a disagreement in overall variabilities. This may be due to the disagreement in either between-item variabilities or within-item variabilities, or both. Roy (2009) allows for a formal test of each.

1.8 Roy's Approach

For the purposes of comparing two methods of measurement, Roy (2009) presents a methodology utilizing linear mixed effects model. This methodology provides for the formal testing of inter-method bias, between-subject variability and within-subject variability of two methods.

Roy (2009) uses an approach based on linear mixed effects (LME) models for the purpose of comparing the agreement between two methods of measurement, where replicate measurements on items (often individuals) by both methods are available. Three tests of hypothesis appropriate are provided for evaluating the agreement between the two methods of measurement under this sampling scheme.

Importantly Roy (2009) further proposes a series of three tests on the variance components of an LME model, which allow decisions on the second and third of Roy's criteria. For these tests, four candidate LME models are constructed.

Roy (2009) proposes the use of LME models to perform a test on two methods of agreement to comparing the agreement between two methods of measurement, where replicate measurements on items (often individuals) by both methods are available, determining whether they can be used interchangeably. The methodology proposed by Roy (2009) is largely based on Hamlett et al. (2004), which in turn follows on from Lam et al. (1999).

Roy (2009) proposes a novel method using the LME model with Kronecker product covariance struc-

ture in a doubly multivariate set-up to assess the agreement between a new method and an established method with unbalanced data and with unequal replications for different subjects. By doubly multivariate set up, Roy means that the information on each patient or item is multivariate at two levels, the number of methods and number of replicated measurements. Further to Lam et al. (1999), it is assumed that the replicates are linked over time. However it is easy to modify to the unlinked case.

The well-known “Limits of Agreement”, as developed by Bland and Altman (1986) are not referred to directly, but are easily computable using the framework proposed by Roy (2009). Further discussion will be provided in due course.

Further to this, Roy (2009) demonstrates an suite of tests that can be used to determine how well two methods of measurement, in the presence of repeated measures, agree with each other.

- No Significant inter-method bias
- No difference in the between-subject variabilities of the two methods
- No difference in the within-subject variabilities of the two methods

The formulation presented above usefully facilitates a series of significance tests that advise as to how well the two methods agree. These tests are as follows:

- A formal test for the equality of between-item variances,
- A formal test for the equality of within-item variances,
- A formal test for the equality of overall variances.

These tests are complemented by the ability to consider the inter-method bias and the overall correlation coefficient. Two methods can be considered to be in agreement if criteria based upon these methodologies are met. Additionally Roy makes reference to the overall correlation coefficient of the two methods, which is determinable from variance estimates. Other important aspects of the method comparison study are consequent. The limits of agreement are computed using the results of the reference model.

Differences in between-subject variabilities of the two methods arise when one method is yielding average response levels for individuals than are more variable than the average response levels for the

same sample of individuals taken by the other method. Differences in within-subject variabilities of the two methods arise when one method is yielding responses for an individual than are more variable than the responses for this same individual taken by the other method. The two methods of measurement can be considered to agree, and subsequently can be used interchangeably, if all three null hypotheses are true.

1.8.1 LME Model Specification

Let y_{mir} denote the r th replicate measurement on the i th item by the m th method, where $m = 1, 2$, $i = 1, \dots, N$, and $r = 1, \dots, n_i$. When the design is balanced and there is no ambiguity we can set $n_i = n$. The LME model underpinning Roy's approach can be written

$$y_{mir} = \beta_0 + \beta_m + b_{mi} + \epsilon_{mir}. \quad (1.1)$$

Here β_0 and β_m are fixed-effect terms representing, respectively, a model intercept and an overall effect for method m .

The b_{1i} and b_{2i} terms represent random effect parameters corresponding to the two methods, having $E(b_{mi}) = 0$ with $\text{Var}(b_{mi}) = g_m^2$ and $\text{Cov}(b_{mi}, b_{m'i}) = d_{12}$. The random error term for each response is denoted ϵ_{mir} having $E(\epsilon_{mir}) = 0$, $\text{Var}(\epsilon_{mir}) = \sigma_m^2$, $\text{Cov}(b_{mir}, b_{m'ir}) = \sigma_{12}$, $\text{Cov}(\epsilon_{mir}, \epsilon_{mir'}) = 0$ and $\text{Cov}(\epsilon_{mir}, \epsilon_{m'ir'}) = 0$.

When two methods of measurement are in agreement, there is no significant differences between β_1 and β_2 , g_1^2 and g_2^2 , and σ_1^2 and σ_2^2 . Here β_0 and β_m are fixed-effect terms representing, respectively, a model intercept and an overall effect for method m .

The model can be reparameterized by gathering the β terms together into (fixed effect) intercept terms $\alpha_m = \beta_0 + \beta_m$. The b_{1i} and b_{2i} terms are correlated random effect parameters having $E(b_{mi}) = 0$ with $\text{Var}(b_{mi}) = g_m^2$ and $\text{Cov}(b_{1i}, b_{2i}) = d_{12}$.

The random error term for each response is denoted ϵ_{mir} having $E(\epsilon_{mir}) = 0$, $\text{Var}(\epsilon_{mir}) = \sigma_m^2$, $\text{Cov}(\epsilon_{1ir}, \epsilon_{2ir}) = \sigma_{12}$, $\text{Cov}(\epsilon_{mir}, \epsilon_{mir'}) = 0$ and $\text{Cov}(\epsilon_{1ir}, \epsilon_{2ir'}) = 0$. Two methods of measurement are in complete agreement if the null hypotheses $H_1: \alpha_1 = \alpha_2$ and $H_2: \sigma_1^2 = \sigma_2^2$ and $H_3: g_1^2 = g_2^2$ hold simultaneously. Roy (2009) uses a Bonferroni correction to control the familywise error rate for tests of $\{H_1, H_2, H_3\}$ and account for difficulties arising due to multiple testing.

1.8.2 Variance Covariance Matrices

Under Roy's model, random effects are defined using a bivariate normal distribution. Consequently, the variance-covariance structures can be described using 2×2 matrices. A discussion of the various structures a variance-covariance matrix can be specified under is required before progressing. The following structures are relevant: the identity structure, the compound symmetric structure and the symmetric structure.

The differences in the models are specifically in how the D and Λ matrices are constructed, using either an unstructured form or a compound symmetry form. To illustrate these differences, consider a generic matrix A ,

$$\mathbf{A} = \begin{pmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{pmatrix}.$$

A symmetric matrix allows the diagonal terms a_{11} and a_{22} to differ. The compound symmetry structure requires that both of these terms be equal, i.e $a_{11} = a_{22}$.

The identity structure is simply an abstraction of the identity matrix. The compound symmetric structure and symmetric structure can be described with reference to the following matrix (here in the context of the overall covariance Block- $\mathbf{\Omega}_i$, but equally applicable to the component variabilities \mathbf{D} and $\mathbf{\Sigma}$);

$$\begin{pmatrix} \omega_1^2 & \omega_{12} \\ \omega_{12} & \omega_2^2 \end{pmatrix}$$

Symmetric structure requires the equality of all the diagonal terms, hence $\omega_1^2 = \omega_2^2$. Conversely compound symmetry make no such constraint on the diagonal elements. Under the identity structure, $\omega_{12} = 0$. A comparison of a model fitted using symmetric structure with that of a model fitted using the compound symmetric structure is equivalent to a test of the equality of variance.

Independence

As though analyzed using between subjects analysis.

$$\begin{pmatrix} \psi^2 & 0 & 0 \\ 0 & \psi^2 & 0 \\ 0 & 0 & \psi^2 \end{pmatrix}$$

Compound Symmetry

Assumes that the variance-covariance structure has a single variance (represented by ψ^2) for all 3 of the time points and a single covariance (represented by ψ_{ij}) for each of the pairs of trials.

$$\begin{pmatrix} \psi^2 & \psi_{12} & \psi_{13} \\ \psi_{21} & \psi^2 & \psi_{23} \\ \psi_{31} & \psi_{32} & \psi^2 \end{pmatrix}$$

1.8.3 Model Terms (Roy 2009)

- Let y_{mir} be the response of method m on the i th subject at the r —th replicate.
- Let \mathbf{y}_{ir} be the 2×1 vector of measurements corresponding to the i —th subject at the r —th replicate.
- Let \mathbf{y}_i be the $R_i \times 1$ vector of measurements corresponding to the i —th subject, where R_i is number of replicate measurements taken on item i .
- Let $\alpha_m i$ be the fixed effect parameter for method for subject i .
- Formally ARoy2009 uses a separate fixed effect parameter to describe the true value μ_i , but later combines it with the other fixed effects when implementing the model.
- Let u_{1i} and u_{2i} be the random effects corresponding to methods for item i .
- $\boldsymbol{\epsilon}_i$ is a n_i -dimensional vector comprised of residual components. For the blood pressure data $n_i = 85$.

- β is the solutions of the means of the two methods. In the LME output, the bias and corresponding t-value and p-values are presented. This is relevant to ARoy2009's first test.

1.8.4 Model Terms (Roy 2009)

It is important to note the following characteristics of this model.

Let the number of replicate measurements on each item i for both methods be n_i , hence $2 \times n_i$ responses. However, it is assumed that there may be a different number of replicates made for different items. Let the maximum number of replicates be p . An item will have up to $2p$ measurements, i.e. $\max(n_i) = 2p$.

Later on \mathbf{X}_i will be reduced to a 2×1 matrix, to allow estimation of terms. This is due to a shortage of rank. The fixed effects vector can be modified accordingly.

\mathbf{Z}_i is the $2n_i \times 2$ model matrix for the random effects for measurement methods on item i .

\mathbf{b}_i is the 2×1 vector of random-effect coefficients on item i , one for each method.

ϵ is the $2n_i \times 1$ vector of residuals for measurements on item i .

\mathbf{G} is the 2×2 covariance matrix for the random effects.

\mathbf{R}_i is the $2n_i \times 2n_i$ covariance matrix for the residuals on item i .

The expected value is given as $E(\mathbf{y}_i) = \mathbf{X}_i\beta$. (Hamlett et al., 2004)

The variance of the response vector is given by $\text{Var}(\mathbf{y}_i) = \mathbf{Z}_i\mathbf{G}\mathbf{Z}_i' + \mathbf{R}_i$ (Hamlett et al., 2004).

\mathbf{b}_i is a m -dimensional vector comprised of the random effects.

$$\mathbf{b}_i = \begin{pmatrix} b_{1i} \\ b_{2i} \end{pmatrix} \quad (1.2)$$

\mathbf{V} represents the correlation matrix of the replicated measurements on a given method. Σ is the within-subject VC matrix.

\mathbf{V} and $\mathbf{\Sigma}$ are positive definite matrices. The dimensions of \mathbf{V} and $\mathbf{\Sigma}$ are $3 \times 3 (= p \times p)$ and $2 \times 2 (= k \times k)$.

It is assumed that \mathbf{V} is the same for both methods and $\mathbf{\Sigma}$ is the same for all replications.

$\mathbf{V} \otimes \mathbf{\Sigma}$ creates a $6 \times 6 (= kp \times kp)$ matrix. \mathbf{R}_i is a sub-matrix of this.

Chapter 2

Linear Mixed effects Models

2.1 Case Deletion Diagnostics for LME models

? describes three propositions that are required for efficient case-deletion in LME models. The first proposition describes how to efficiently update V when the i th element is deleted.

$$V_{[i]}^{-1} = \Lambda_{[i]} - \frac{\lambda\lambda'}{\nu_{ii}} \quad (2.1)$$

The second of Christensen's propositions is the following set of equations, which are variants of the Sherman Woodbury updating formula.

$$X'_{[i]} V_{[i]}^{-1} X_{[i]} = X' V^{-1} X - \frac{\hat{x}_i \hat{x}_i'}{s_i} \quad (2.2)$$

$$(X'_{[i]} V_{[i]}^{-1} X_{[i]})^{-1} = (X' V^{-1} X)^{-1} + \frac{(X' V^{-1} X)^{-1} \hat{x}_i \hat{x}_i' (X' V^{-1} X)^{-1}}{s_i - \bar{h}_i} \quad (2.3)$$

$$X'_{[i]} V_{[i]}^{-1} Y_{[i]} = X' V^{-1} Y - \frac{\hat{x}_i \hat{y}_i'}{s_i} \quad (2.4)$$

Influence on measure component ratios

The general diagnostic tools for variance component ratios are the analogues of the Cook's distance and the Information Ratio.

$$\begin{aligned}
CD_U(\gamma) &= (\hat{\gamma}_{(U)} - \hat{\gamma})' [\text{var}(\hat{\gamma})]^{-1} (\hat{\gamma}_{(U)} - \hat{\gamma}) \\
&= -\mathbf{g}'_{(U)} (\mathbf{Q} - \mathbf{G})^{-1} \mathbf{Q} (\mathbf{Q} - \mathbf{G}) \mathbf{g}_{(U)} \\
&= \mathbf{g}'_{(U)} (\mathbf{I}_r \text{var}(\hat{\gamma}) \mathbf{G})^{-2} \text{var}(\hat{\gamma}) \mathbf{g}_{(U)}
\end{aligned}$$

Large values of $CD(\gamma)$ highlight observation groups for closer attentions

$$IR\gamma = \frac{\det(\mathbf{Q} - \mathbf{G})}{\det(\mathbf{Q})}$$

Ideally when all observations have the same influence on the information matrix $IR\gamma$ is approximately one. Deviations from one indicate the group U is influential. Since $\text{var}(\hat{\gamma})$ and \mathbf{I}_r are fixed for all observations, $IR\gamma$ is a function of \mathbf{G} , in turn a function of \mathbf{C}_i and c_{ii} .

2.2 BXC - Model Terms

- Let y_{mir} be the response of method m on the i th subject at the r -th replicate.
- Let \mathbf{y}_{ir} be the 2×1 vector of measurements corresponding to the i -th subject at the r -th replicate.
- Let \mathbf{y}_i be the $R_i \times 1$ vector of measurements corresponding to the i -th subject, where R_i is number of replicate measurements taken on item i .
- Let $\alpha_m i$ be the fixed effect parameter for method for subject i .
- Formally Roy uses a separate fixed effect parameter to describe the true value μ_i , but later combines it with the other fixed effects when implementing the model.
- Let u_{1i} and u_{2i} be the random effects corresponding to methods for item i .
- ϵ_i is a n_i -dimensional vector comprised of residual components. For the blood pressure data $n_i = 85$.

- β is the solutions of the means of the two methods. In the LME output, the bias and corresponding t-value and p-values are presented. This is relevant to Roy's first test.

2.3 LME

Consistent with the conventions of mixed models, ? formulates the measurement y_{ij} from method i on individual j as follows;

$$y_{ij} = P_{ij}\theta + W_{ij}v_i + X_{ij}b_j + Z_{ij}u_j + \epsilon_{ij}, (j = 1, 2, i = 1, 2, \dots, n) \quad (2.5)$$

The design matrix P_{ij} , with its associated column vector θ , specifies the fixed effects common to both methods. The fixed effect specific to the j th method is articulated by the design matrix W_{ij} and its column vector v_i . The random effects common to both methods is specified in the design matrix X_{ij} , with vector b_j whereas the random effects specific to the i th subject by the j th method is expressed by Z_{ij} , and vector u_j . Noticeably this notation is not consistent with that described previously. The design matrices are specified so as to include a fixed intercept for each method, and a random intercept for each individual. Additional assumptions must also be specified;

$$v_{ij} \sim N(0, \Sigma), \quad (2.6)$$

These vectors are assumed to be independent for different i s, and are also mutually independent. All Covariance matrices are positive definite. In the above model effects can be classed as those common to both methods, and those that vary with method. When considering differences, the effects common to both effectively cancel each other out. The differences of each pair of measurements can be specified as following;

$$d_{ij} = X_{ij}b_j + Z_{ij}u_j + \epsilon_{ij}, (j = 1, 2, i = 1, 2, \dots, n) \quad (2.7)$$

This formulation has separate distributional assumption from the model stated previously.

This agreement covariate x is the key step in how this methodology assesses agreement.

2.4 Remarks

The relationship between precision and the within-item and between-item variability must be established. Roy establishes the equivalence of repeatability and within-item variability, and hence precision. The method with the smaller within-item variability can be deemed to be the more precise.

A useful approach is to compute the confidence intervals for the ratio of within-item standard deviations (equivalent to the ratio of repeatability coefficients), which can be interpreted in the usual manner.

In fact, the ratio of within-item standard deviations, with the attendant confidence interval, can be determined using a single R command: `intervals()`.

Pinheiro and Bates (pg 93-95) give a description of how confidence intervals for the variance components are computed. Furthermore a complete set of confidence intervals can be computed to complement the variance component estimates.

What is required is the computation of the variance ratios of within-item and between-item standard deviations.

A naive approach would be to compute the variance ratios by relevant F distribution quantiles. However, the question arises as to the appropriate degrees of freedom. Limits of agreement are easily computable using the LME framework. While we will not be considering this analysis, a demonstration will be provided in the example.

Chapter 3

LME Likelihood

3.1 PRESS

An (unconditional) predicted value is $\hat{y}_i = x_i' \hat{\beta}$, where the vector x_i is the i th row of \mathbf{X} .

An (unconditional) predicted value is $\hat{y}_i = x_i' \hat{\beta}$, where the vector x_i is the i th row of \mathbf{X} . The (raw) residual is given as $\varepsilon_i = y_i - \hat{y}_i$. The PRESS residual is similarly constructed, using the predicted value for observation i with a model fitted from reduced data.

$$\varepsilon_{i(U)} = y_i - x_i' \hat{\beta}_{(U)}$$

3.2 One Way ANOVA

3.2.1 Page 448

Computing the variance of $\hat{\beta}$

$$\text{var}(\hat{\beta}) = (X'V^{-1}X)^{-1} \quad (3.1)$$

It is not necessary to compute V^{-1} explicitly.

$$V^{-1}X = \Sigma^{-1}X - Z(Z'\Sigma^{-1}X) \quad (3.2)$$

$$= \Sigma^{-1}(X - Zb_x) \quad (3.3)$$

The estimate b_x is the same term obtained from the random effects model; $X = Zb_x + e$, using X as an outcome variable. This formula is convenient in applications where b_x can be easily computed. Since X is a matrix of p columns, b_x can simple be computed columnn by column. according to the columns of X .

3.2.2 Page 448- simple example

Consider a simple model of the form;

$$y_{ij} = \mu + \beta_i + \epsilon_{ij}.$$

The iterative procedure is as follows Evaluate the individual group mean \bar{y}_i and variance \hat{Sigma}_i^2 . Then use the variance of the group means as an estimate of the σ_b^2 . The average of the the variances of the groups is the initial estimate of the σ_e^2 .

Iterative procedure

The iterative procedure comprises two steps, with 0 as the first approximation of b_i .

The first step is to compute λ , the ratio of variabilities,

$$\lambda = \frac{\sigma_b^2}{\sigma_e^2}$$

$$\begin{aligned}\mu &= \frac{1}{N} \sum_{ij} (y_{ij} - b_i) \\ b_i &= \frac{n(\bar{y}_i - \mu)}{n + \lambda}\end{aligned}$$

The second step is to updat σ_e^2

$$\sigma_e^2 = \frac{e'e}{N - df} \tag{3.4}$$

where e is the vector of $e_{ij} = y_{ij} - \mu - b_i$ and $df = qn/n + \lambda$ and

$$\sigma_b^2 = \frac{1}{q} \sum_{i=1}^q b_i^2 + \left(\frac{n}{\sigma_e^2} + \frac{1}{\sigma_b^2} \right)^{-1} \tag{3.5}$$

Worked Example

Further to [pawitan 17.1] the initial estimates for variability are $\sigma_b^2 = 1.7698$ and $\sigma_e^2 = 0.3254$. At convergence the following results are obtained.

n=16, q=5

$$\hat{\mu} = \bar{y} = 14.175$$

$$\hat{\sigma}^2 = 0.325$$

$$\hat{\sigma}_b^2 = 1.395$$

$$\sigma = 0.986$$

At convergence the following estimates are obtained,

$$\hat{\mu} = 14.1751$$

$$\hat{b} = (-0.6211, 0.2683, 1.4389, -1.914, 0.8279)$$

$$\hat{\sigma}_b^2 = 1.3955$$

$$\hat{\sigma}_e^2 = 0.3254$$

3.3 Sampling

One important feature of replicate observations is that they should be independent of each other. In essence, this is achieved by ensuring that the observer makes each measurement independent of knowledge of the previous value(s). This may be difficult to achieve in practice. (Check who said this)

3.4 Conclusion

Carstensen et al. (2008) and Roy (2009) highlight the need for method comparison methodologies suitable for use in the presence of replicate measurements. Roy (2009) presents a comprehensive methodology for assessing the agreement of two methods, for replicate measurements. This methodology has

the added benefit of overcoming the problems of unbalanced data and unequal numbers of replicates. Implementation of the methodology, and interpretation of the results, is relatively easy for practitioners who have only basic statistical training. Furthermore, it can be shown that widely used existing methodologies, such as the limits of agreement, can be incorporated into ARoy2009’s methodology.

Permutation Test, Power Tests and Missing Data

This section explores topics such as dependent variable simulation and power analysis, introduced by Galecki & Burzykowski (2013), and implementable with their *nlmeU* R package.

Using the *predictmeans* R package, it is possible to perform permutation t-tests for coefficients of (fixed) effects and permutation F-tests.

The matter of missing data has not been commonly encountered in either Method Comparison Studies or Linear Mixed Effects Modelling. However ARoy2009 (2009) deals with the relevant assumptions regarding missing data.

Galecki & Burzykowski (2013) approaches the subject of missing data in LME Modelling. The *nlmeU* package includes the `patMiss` function, which “*allows to compactly present pattern of missing data in a given vector/matrix/data frame or combination of thereof*”.

3.4.1 EBLUPS-Diagnostics for Random Effects

West et al. (2007) recommends the empirical Bayes predictor, also known as EBLUPS as a diagnostic tool for Random effects. Checking EBLUPS for normality is of limited value.

The EBLUP is useful to identify outlier subjects given that it represents the distance between the population mean value and the value predicted for the i th subject. A way of using the EBLUP to search for outliers subjects is to use the Mahalanobis distance (see Waternaux et al., 1989), FORMULA

. It is also possible to use the EBLUP to verify the random effects normality assumption. For more information; see Nobre and Singer (2007). In Table 2 we summarize diagnostic techniques involving residuals discussed in Nobre and Singer (2007).

Chapter 4

General Appendices

$$\Lambda = \frac{\max_{H_0} L}{\max_{H_1} L}$$

4.0.1 Extending deletion diagnostics to LMEs

? notes the case deletion diagnostics techniques have not been applied to linear mixed effects models and seeks to develop methodologies in that respect. ? develops these techniques in the context of REML

$$X = \begin{bmatrix} x'_i \\ X(i) \end{bmatrix}, Z = \begin{bmatrix} z'_{ij} \\ Z_{j(i)} \end{bmatrix}, Z = \begin{bmatrix} z'_{ij} \\ Z_{j(i)} \end{bmatrix},$$

$$y = \begin{bmatrix} y'_{ij} \\ y_{j(i)} \end{bmatrix} \text{ and } H = \begin{bmatrix} h_{ii} & h \\ h_{j(i)} & h \end{bmatrix}$$

For notational simplicity, $\mathbf{A}_{(i)}$ denotes an $n \times m$ matrix \mathbf{A} with the i -th row removed, \mathbf{a}_i denotes the i -th row of \mathbf{A} , and a_{ij} denotes the (i, j) -th element of \mathbf{A} .

$\mathbf{a}_{(i)}$ denotes a vector \mathbf{a} with the i -th element, a_i , removed.

$$\check{\mathbf{a}}_i = \mathbf{a}_i - \mathbf{A}_{(i)} \mathbf{H}_{[i]} \mathbf{h}_i \quad (4.1)$$

4.1 Unknown Material

To standardize the assessment of how influential data is, several measures of influence are commonly used, such as DFBETAS and Cooks Distance.

Although influential cases thus have extreme values on one or more of the variables, they can be onliers rather than outliers.

To account for this, the (standardized) deleted residual is defined as the difference between the observed score of a case on the dependent variable, and the predicted score from the regression model fitted from data when that case is omitted.

Just as influential cases are not necessarily outliers, outliers are not necessarily influential cases.

This also holds for deleted residuals. The reason for this is that the amount of influence a case exerts on the regression slope is not only determined by how well its (observed) score is fitted by the specified regression model, but also by its score(s) on the independent variable(s). The degree to which the scores of a case on the independent variable(s) are extreme is indicated by the leverage of this case.

4.1.1 Estimation

$$\hat{\beta} = X^T \quad (4.2)$$

$$\hat{\gamma} = G(\hat{\theta})Z^T \quad (4.3)$$

The difference between perturbation and residual analysis between the linear and LME models. The estimates of the fixed effects β depend on the estimates of the covariance parameters.

4.1.2 Zewotir-Cook's Distance

Diagnostic tool for variance components

$$C_{\theta i} = ((\hat{\theta})_{[i]} - \hat{\theta})^T \text{cov}(\hat{\theta})^{-1} ((\hat{\theta})_{[i]} - \hat{\theta})$$

linear functions $CD(\psi)_i$ does not have to be calculated unless $CD(\beta)_i$ is large.

It is also desirable to measure the influence of the case deletions on the covariance matrix of $\hat{\beta}$.

For fixed effects parameter estimates in LME models, the Cook's distance can be extended to measure influence on these fixed effects.

$$CD_i(\beta) = \frac{(c_{ii} - r_{ii}) \times t_i^2}{r_{ii} \times p}$$

For random effect estimates, the Cook's distance is

$$CD_i(b) = g'_{(i)}(I_r + \text{var}(\hat{b})D)^{-2} \text{var}(\hat{b})g_{(i)}.$$

Large values for Cook's distance indicate observations for special attention.

For linear functions, $CD(\psi)_i$ does not have to be calculated unless $CD(\beta)_i$ is large.

Mean Square Prediction Error

$$MSPR = \frac{\sum (y_i - \hat{y}_i)^2}{n^*} \quad (4.4)$$

4.1.3 Leverage

Leverage can be defined through the projection matrix that results from a transformation of the model with the inverse of the Cholesky decomposition of \mathbf{V} , or an oblique projector: $\mathbf{Y} = \mathbf{H}\hat{\mathbf{Y}}$.

While H is idempotent, it is generally not symmetric and thus not a projection matrix in the narrow sense.

$$h_{ii} = x_i'(X'X)^{-1}x_i$$

The trace of \mathbf{H} equals the rank of \mathbf{X} . If V_{ij} denotes the element in row i , column j of \mathbf{V}^{-1} , then for a model containing only an intercept the diagonal elements of \mathbf{H} .

$$h_{ii} = \frac{\sum v_{ij}}{\sum \sum v_{ij}}$$

PRESS

Schabenberger (2004) describes the use of the *PRESS* and *DFITS* in determining influence.

The *PRESS* residual is the difference between the observed value and the predicted (marginal) value.

$$\hat{e}_{i(U)} = y_i - x_i\hat{\beta}_{(U)} \quad (4.5)$$

The prediction residual sum of squares (PRESS) is a value associated with this calculation. When fitting linear models, PRESS can be used as a criterion for model selection, with smaller values indicating better model fits.

$$PRESS = \sum (y - y^{(k)})^2 \quad (4.6)$$

$$e_{-Q} = y_Q - x_Q\hat{\beta}^{-Q}$$

$$PRESS = \sum (y - y^{-Q})^2$$

$$PRESS_{(U)} = y_i - x_i\hat{\beta}_{(U)}$$

PRESS Residuals and PRESS Statistic

The predicted residual sum of squares (PRESS) statistic is a form of cross-validation used in regression analysis to provide a summary measure of the fit of a model to a sample of observations that were

not themselves used to estimate the model. It is calculated as the sums of squares of the prediction residuals for those observations.

A fitted model having been produced, each observation in turn is removed and the model is refitted using the remaining observations. The out-of-sample predicted value is calculated for the omitted observation in each case, and the PRESS statistic is calculated as the sum of the squares of all the resulting prediction errors:[4]

$$\text{PRESS} = \sum_{i=1}^n (y_i - \hat{y}_{i,-i})^2$$

Given this procedure, the PRESS statistic can be calculated for a number of candidate model structures for the same dataset, with the lowest values of PRESS indicating the best structures. Models that are over-parameterised (over-fitted) would tend to give small residuals for observations included in the model-fitting but large residuals for observations that are excluded.

4.1.4 Local Influence

? developed their global influences for the deletion of single observations in two steps: a one-step estimate for the REML (or ML) estimate of the variance components, and an ordinary case-deletion diagnostic for a weighted regression problem (conditional on the estimated covariance matrix) for fixed effects.

Bibliography

- ACR (2008). Acute Chest Pain (suspected aortic dissection) - American College of Radiology Expert Group Report.
- Akaike, H. (1974). A new look at the statistical model identification. *IEEE Transactions on Automatic Control* 19(6), 716–723.
- Altman, D. and J. Bland (1983). Measurement in medicine: The analysis of method comparison studies. *Journal of the Royal Statistical Society. Series D (The Statistician)* 32(3), 307–317.
- Barnhart, H., M. Haber, and L. Lin (2007). An overview of assessing agreement with continuous measurements. *Journal of Biopharmaceutical Statistics* 17, 529–569.
- Bartko, J. (1994). Measures of agreement: A single procedure. *Statistics in Medicine* 13, 737–745.
- Benjamini, Y. and Y. Hochberg (1995). Controlling the false discovery rate: A practical and powerful approach to multiple testing. *J. Roy. Statist. Soc. Ser. B* 57(1), 289–300.
- Bland, J. and D. Altman (1986). Statistical methods for assessing agreement between two methods of clinical measurement. *The Lancet* i, 307–310.
- Bland, J. and D. Altman (1995). Comparing methods of measurement - why plotting difference against standard method is misleading. *The Lancet* 346, 1085–87.
- Bland, J. and D. Altman (1999). Measuring agreement in method comparison studies. *Statistical Methods in Medical Research* 8(2), 135–160.
- Brown, H. and R. Prescott (1999). *Applied Mixed Models In Medicine*. John Wiley and Sons.

- Carstensen, B. (2004). Comparing and predicting between several methods of measurement. *Biostatistics* 5(3), 399–413.
- Carstensen, B. (2010). *Comparing Clinical Measurement Methods : A Practical Guide*. Wiley.
- Carstensen, B., J. Simpson, and L. C. Gurrin (2008). Statistical models for assessing agreement in method comparison studies with replicate measurements. *The International Journal of Biostatistics* 4(1).
- Cook, R. (1977). Detection of influential observations in linear regression. *Technometrics* 19, 15–18.
- Cornbleet, P. J. and D. Cochrane (1979). Regression methods for assessing agreement between two methods of clinical measurement. *Journal of Clinical Chemistry* 24(2), 342–345.
- Demidenko, E. (2004). *Mixed Models: Theory And Application*. Dartmouth College: Wiley Interscience.
- Dempster, A., N. Laird, and D. Rubin (1977). Maximum likelihood from incomplete data via the em algorithm. *ournal of the Royal Statistical Society. Series B* 39(1), 1–38.
- Dewitte, K., C. Fierens, D. Stckl, and L. M. Thienpont (2002). Application of the Bland Altman plot for interpretation of method-comparison studies: A critical investigation of its practice. *Clinical Chemistry* 48, 799–801.
- Dodge, Y. (2003). *The Oxford Dictionary of Statistical Terms*. Oxford University Press.
- Dunn, G. (2002). *Statistical Evaluation of Measurement Error* (Second ed.). Stanford: American Mathematical Society and Digital Press.
- Eisenhart, C. (1947). The assumptions underlying the analysis of variance. *Biometrics* 3(1), 1–21.
- Fisher, R. (1918). The correlation between relatives on the supposition of mendelian inheritance. *Transactions of the Royal Society of Edinburgh* 2, 399–433.
- Grubbs, F. (1973). Errors of measurement, precision, accuracy and the statistical comparison of measuring instruments. *Technometrics* 15(1), 53–66.

- Hamlett, A., L. Ryan, and R. Wolfinger (2004). On the use of PROC MIXED to estimate correlation in the presence of repeated measures. *Proceedings of the Statistics and Data Analysis Section, SAS Users Group International 198-229*, 1–7.
- Hartley, H. and J. Rao (1967). Maximum-likelihood estimation for the mixed analysis of variance model. *Biometrika* 54(1/2), 93–108.
- Harville, D. A. (1977). Maximum likelihood approaches to variance component estimation and to related problems. *Journal of American Statistical Association* 72(358), 320–338.
- Haslett, J. (1999). A simple derivation of deletion diagnostic results for the general linear model with correlated errors. *Journal of the Royal Statistical Society (Series B)* 61, 603–609.
- Haslett, J. and D. Dillane (2004). Application of ‘delete = replace’ to deletion diagnostics for variance component estimation. *Journal of the Royal Statistical Society (Series B)* 66, 131–143.
- Haslett, J. and K. Hayes (1998). Residuals for the linear model with general covariance structure. *Journal of the Royal Statistical Society (Series B)* 60, 201–215.
- Henderson, C. (1953). Estimation of variance and covariance components. *Biometrics* 9(2), 226–252.
- Henderson, C., O. Kempthorne, S. Searle, and C. von Krosigk (1959). The estimation of environmental and genetic trends from records subject to culling. *Biometrics* 15, 192–218.
- Henderson, C., O. Kempthorne, S. Searle, and C. von Krosigk (1963). The estimation of environmental and genetic trends from records subject to culling. *Biometrics* 15, 192–218.
- Henderson, C., O. Kempthorne, S. Searle, and C. von Krosigk (1973). The estimation of environmental and genetic trends from records subject to culling. *Biometrics* 15, 192–218.
- Henderson, C., O. Kempthorne, S. Searle, and C. von Krosigk (1984). The estimation of environmental and genetic trends from records subject to culling. *Biometrics* 15, 192–218.
- Henderson, C. R. (1950). Estimation of genetic parameters (abstract). *Annals of Mathematical Statistics* 21, 309–310.

- Laird, N. M. and J. H. Ware (1982). Random-effects models for longitudinal data. *Biometrics* 38(4), 963–974.
- Lam, M., K. Webb, and D. O'Donnell (1999). Correlation between two variables in repeated measurements. *American Statistical Association, Proceedings of the Biometric Session*, 213–218.
- Lewis, P., P. Jones, J. Polak, and H. Tillitson (1991). The problem of conversion in method comparison studies. *Applied Statistics* 40, 105–112.
- Lin, L. (2000). Total deviation index for measuring individual agreement with applications in laboratory performance and bioequivalence. *Statistics in medicine* 19, 255–270.
- Lin, L., A. Hedayat, B. Sinha, and M. Yang (2002). Statistical methods in assessing agreement: Models issues and tools. *Journal of American Statistical Association* 97, 257–270.
- Lin, S. C., D. M. Whipple, and Charles S Ho (1991). Evaluation of statistical equivalence using limits of agreement and associates sample size calculation. *Communications in Statistics - Theory and Methods* 20(6), 1419–1432.
- Linnet, K. (1998). Performance of deming regression analysis in case of misspecified analytical error ratio in method comparison studies. *Clinical Chemistry* 44, 1024–1031.
- Ludbrook, J. (1997). Comparing methods of measurement. *Clinical and Experimental Pharmacology and Physiology* 24, 193–203.
- Ludbrook, J. (2002). Statistical techniques for comparing measurers and methods of measurement: a critical review. *Clinical and Experimental Pharmacology and Physiology* 29, 527–536.
- Luiz, R., A. Costa, P. Kale, and G. Werneck (2003). Assessment of agreement of a quantitative variable: a new graphical approach. *Journal of Clinical Epidemiology* 56, 963–967.
- Mantha, S., M. F. Roizen, L. A. Fleisher, R. Thisted, and J. Foss (2000). Comparing methods of clinical measurement: Reporting standards for bland and altman analysis. *Anaesthesia and Analgesia* 90, 593–602.

- McCullough, C. and S. Searle (2001). *Generalized , Linear and Mixed Models*. Wiley Interscience.
- Myles, P. (2007). Using the Bland Altman method to measure agreement with repeated measures. *British Journal of Anaesthesia* 99(3), 309–311.
- NIST (2009). Cesium fountain atomic clock : The primary time and frequency standard for the United States. <http://tf.nist.gov/timefreq/cesium/fountain.htm>.
- Paterson, H. and R. Thompson (1971). Recovery of inter-block information when block sizes are unequal. *Biometrika* 58(3), 545–554.
- Patterson, H. and R. Thompson (1971). Recovery of inter-block information when block sizes are unequal. *Biometrika* 58(3), 545–554.
- Phelps, C. and A. Hutson (1995). Estimating diagnostic test accuracy using a fuzzy gold standard. *Medical decision making* 15, 144–57.
- Pinheiro, J. and D. Bates (1994). *Mixed Effects Models in S and S plus* (2nd ed.). Reading, Massachusetts: Springer.
- Pizzi, N. (1999). Fuzzy pre-processing of gold standards as applied to biomedical spectra classification. *Artificial Intelligence in Medicine* 16, 171–182.
- Roy, A. (2009). An application of the linear mixed effects model to ass the agreement between two methods with replicated observations. *Journal of Biopharmaceutical Statistics* 19, 150–173.
- Schabenberger, O. (2004). Mixed model influence diagnostics. 18929.
- Searle, S. (1997). *Linear Models*. Wiley classics Library.
- Singer, J. D. (1998). Using sas proc mixed to fit multilevel models, hierarchical models, and individual growth models. *Journal of Educational and Behavioral Statistics* 24(4), 323–355.
- Tippett, L. (1931). *The Methods of Statistics* (1st ed.). London: Williams and Norgate.

West, B., K. Welch, and A. Galecki (2007). *Linear Mixed Models: a Practical Guide Using Statistical Software*. Chapman and Hall CRC.

Zewotir, T. and J. Galpin (2005). Influence diagnostics for linear mixed models. *Journal of Data Science* 3, 153–177.

In the graph above, you can predict non-zero values for the residuals based on the fitted value. For example, a fitted value of 8 has an expected residual that is negative. Conversely, a fitted value of 5 or 11 has an expected residual that is positive.

The non-random pattern in the residuals indicates that the deterministic portion (predictor variables) of the model is not capturing some explanatory information that is leaking into the residuals. The graph could represent several ways in which the model is not explaining all that is possible.

Possibilities include:

- A missing variable
- A missing higher-order term of a variable in the model to explain the curvature
- A missing interaction between terms already in the model

Identifying and fixing the problem so that the predictors now explain the information that they missed before should produce a good-looking set of residuals.

In addition to the above, here are two more specific ways that predictive information can sneak into the residuals:

The residuals should not be correlated with another variable. If you can predict the residuals with another variable, that variable should be included in the model. In Minitabs regression, you can plot the residuals by other variables to look for this problem.

Autocorrelation

Adjacent residuals should not be correlated with each other (**autocorrelation**). If you can use one residual to predict the next residual, there is some predictive information present that is not captured by the predictors. Typically, this situation involves time-ordered observations. For example, if a residual is more likely to be followed by another residual that has the same sign, adjacent residuals are positively

correlated. You can include a variable that captures the relevant time-related information, or use a time series analysis.

In Minitabs regression, you can perform the ***Durbin-Watson*** test to test for autocorrelation.

4.2 ICC, Reproducibility Index and Passing-Bablok

4.2.1 Intraclass Correlation Coefficient

This measure of agreement is estimated using variance components from appropriate analysis of variance models. Measures of agreement are variance dependent, and so the ICC can be misleading. The ICC takes a value between 0 and 1, and is based on Analysis of Variance methodologies.

The ICC is a measure of reliability. Bartko (1994) considers the ICC as just another measure of agreement.

Intra-class correlation coefficient

The ICC, which takes on values between 0 and 1, is based on analysis of variance techniques. It is close to 1 when the differences between paired measurements is very small compared to the differences between subjects. Of these three procedures—t test, correlation coefficient, intra-class correlation coefficient—the ICC is best because it can be large only if there is no bias and the paired measurements are in good agreement, but it suffers from the same faults ii and iii as ordinary correlation coefficients. The magnitude of the ICC can be manipulated by the choice of samples to split and says nothing about the magnitude of the paired differences.

4.2.2 Passing and Bablok (1983)

Passing & Bablok have described a linear regression model that are without the usual assumptions regarding the distribution of the samples and the measurement errors. The result does not depend on the assignment of the methods (or instruments) to X and Y. The slope and intercept are calculated with their 95% confidence interval. Hypothesis tests on the slope and intercept maybe then carried out. If the hypothesis of the intercept is rejected, then it is concluded that it is significant different from 0

and both raters differ at least by a constant amount.

If the hypothesis of the slope is rejected, then it is concluded that the slope is significantly different from 1 and there is at least a proportional difference between the two raters.

4.2.3 Lin's Reproducibility Index

Lin proposes the use of a reproducibility index, called the Concordance Correlation Coefficient (CCC). While it is not strictly a measure of agreement as such, it can form part of an overall method comparison methodology.

4.3 Repeated Measurements

In cases where there are repeated measurements by each of the two methods on the same subjects, Bland Altman suggest calculating the mean for each method on each subject and use these pairs of means to compare the two methods. The estimate of bias will be unaffected using this approach, but the estimate of the standard deviation of the differences will be too small, because of the reduction of the effect of repeated measurement error. Bland Altman propose a correction for this. Carstensen attends to this issue also, adding that another approach would be to treat each repeated measurement separately.

In this model, the variances of the random effects must depend on m , since the different methods do not necessarily measure on the same scale, and different methods naturally must be assumed to have different variances. Carstensen (2004) attends to the issue of comparative variances.

4.4 Linnet - References

The statistical procedures are described in: Linnet K. Necessary sample size for method comparison studies based on regression analysis. Clin Chem 1999; 45: 882-94. Linnet K. Performance of Deming regression analysis in case of misspecified analytical error ratio in method comparison studies. Clin Chem 1998; 44: 1024-1031. Linnet K. Evaluation of regression procedures for methods comparison studies.

Clin Chem 1993; 39. 424-432. Linnet K. Estimation of the linear relationship between measurements of two methods with proportional errors. Stat Med 1990; 9: 1463-1473.

4.5 Lewis Conversion

While regarding a comparison of two pump meters under operational conditions

..It is suspected that the various assumptions made by each method are weak under operational conditions Lewis listed several sources of variation that relate to the practical aspects of each measurement method.

There is little reasons to believe that the laboratory conditions of the devise provide a suitable basis for the conversion of data gathered under operational conditions.

Latent variables are variables that are not measured (i.e. not observed) but whose values is observed from other observed variables. One advantage of using latent variables is that it reduces the dimensionality of data. A large number of observable variables can be aggregated in a model to represent an underlying concept, making it easier for humans to understand the data. [wikipedia]

4.6 RSquared for LME models

As a complement to this, one can also consider how to properly employ the R^2 measure, in the context of Methoc Comparison Studies, further to the work by Edwards et al, namely “An R^2 statistic for fixed effects in the linear mixed model”.

Abstract for “An R^2 statistic for fixed effects in the linear mixed model”

Statisticians most often use the linear mixed model to analyze Gaussian longitudinal data.

The value and familiarity of the R^2 statistic in the linear univariate model naturally creates great interest in extending it to the linear mixed model. We define and describe

how to compute a model R2 statistic for the linear mixed model by using only a single model.

The proposed R2 statistic measures multivariate association between the repeated outcomes and the fixed effects in the linear mixed model. The R2 statistic arises as a function of an appropriate F statistic for testing all fixed effects (except typically the intercept) in a full model.

The statistic compares the full model with a null model with all fixed effects deleted (except typically the intercept) while retaining exactly the same covariance structure.

Furthermore, the R2 statistic leads immediately to a natural definition of a partial R2 statistic. A mixed model in which ethnicity gives a very small p-value as a longitudinal predictor of blood pressure (BP) compellingly illustrates the value of the statistic.

In sharp contrast to the extreme p-value, a very small R^2 , a measure of statistical and scientific importance, indicates that ethnicity has an almost negligible association with the repeated BP outcomes for the study.

4.7 Remarks on the Multivariate Normal Distribution

Diligence is required when considering the models. Carstensen specifies his models in terms of the univariate normal distribution. ARoy2009's model is specified using the bivariate normal distribution. This gives rise to a key difference between the two models, in that a bivariate model accounts for covariance between the variables of interest. The multivariate normal distribution of a k -dimensional random vector $X = [X_1, X_2, \dots, X_k]$ can be written in the following notation:

$$X \sim \mathcal{N}(\mu, \Sigma),$$

or to make it explicitly known that X is k -dimensional,

$$X \sim \mathcal{N}_k(\mu, \Sigma).$$

with k -dimensional mean vector

$$\mu = [E[X_1], E[X_2], \dots, E[X_k]]$$

and $k \times k$ covariance matrix

$$\Sigma = [\text{Cov}[X_i, X_j]], \quad i = 1, 2, \dots, k; \quad j = 1, 2, \dots, k$$

1. Univariate Normal Distribution

$$X \sim \mathcal{N}(\mu, \sigma^2),$$

2. Bivariate Normal Distribution

(a)

$$X \sim \mathcal{N}_2(\mu, \Sigma),$$

(b)

$$\mu = \begin{pmatrix} \mu_x \\ \mu_y \end{pmatrix}, \quad \Sigma = \begin{pmatrix} \sigma_x^2 & \rho\sigma_x\sigma_y \\ \rho\sigma_x\sigma_y & \sigma_y^2 \end{pmatrix}.$$

4.7.1 Lin's Reproducibility Index

Lin proposes the use of a reproducibility index, called the Concordance Correlation Coefficient (CCC). While it is not strictly a measure of agreement as such, it can form part of an overall method comparison methodology.

4.8 Measurement Error Models

Dunn (2002) proposes a measurement error model for use in method comparison studies. Consider n pairs of measurements X_i and Y_i for $i = 1, 2, \dots, n$.

$$X_i = \tau_i + \delta_i \tag{4.7}$$

$$Y_i = \alpha + \beta\tau_i + \epsilon_i$$

In the above formulation is in the form of a linear structural relationship, with τ_i and $\beta\tau_i$ as the true values, and δ_i and ϵ_i as the corresponding measurement errors. In the case where the units of measurement are the same, then $\beta = 1$.

$$E(X_i) = \tau_i \tag{4.8}$$

$$E(Y_i) = \alpha + \beta\tau_i$$

$$E(\delta_i) = E(\epsilon_i) = 0$$

The value α is the inter-method bias between the two methods.

$$z_0 = d = 0 \tag{4.9}$$

$$z_{n+1} = z_n^2 + c \tag{4.10}$$

4.8.1 The Problem of Identifiability

Dunn (2002) highlights an important issue regarding using models such as these, the identifiability problem. This comes as a result of there being too many parameters to be estimated. Therefore assumptions about some parameters, or estimators used, must be made so that others can be estimated.

For example α may take the value of the inter-method bias estimate from Bland - Altman methodology.

For example in literature the variance ratio $\lambda = \frac{\sigma_1^2}{\sigma_2^2}$ must often be assumed to be equal to 1 (Linnet, 1998). Dunn (2002) considers methodologies based on two methods with single measurements on each subject as inadequate for a serious study on the measurement characteristics of the methods. This is because there would not be enough data to allow for a meaningful analysis. There is, however, a contrary argument that in many practical settings it is very difficult to get replicate observations when the measurement method requires invasive medical procedure.

4.8.2 Identifiability

In many models, naive assumptions are required to overcome issues of identifiability. Precision is defined by the reciprocal of the variance of the random errors. Also it is assumed that the error variance is independent of the amount of material being measured. However, in practice, this is often not the case. Variability increases over the scale of measurements over many cases. Estimators of scale parameters are estimable only if the analyst is prepared to make naive, if not unacceptable, assumptions.

Equation 4 ψ and ε are statistically independent of each other. Contamination effect that arises from non-specificity / specimen specific bias. Random error is measured by . Homogeneity of variances is assumed. If there are no replicate measures, both variances are completely confounded, and there is no way of telling them apart. Scaling of new measurements is measured by .

4.9 Carstensen Model (mir model)

A measurement y_{mi} by method m on individual i is formulated as follows;

$$y_{mi} = \alpha_m + \mu_i + e_{mi} \quad e_{mi} \sim \mathcal{N}(0, \sigma_m^2) \quad (4.11)$$

The differences are expressed as $d_i = y_{1i} - y_{2i}$. For the replicate case, an interaction term c is added to the model, with an associated variance component. All the random effects are assumed independent, and that all replicate measurements are assumed to be exchangeable within each method.

The following model (in the authors own notation) is formulated as follows, where y_{mir} is the r th replicate measurement on subject i with method m .

Using Carstensen's notation, a measurement y_{mi} by method m on individual i the measurement y_{mir} is the r th replicate measurement on the i th item by the m th method, where $m = 1, 2$, $i = 1, \dots, N$, and $r = 1, \dots, n_i$ is formulated as follows;

$$y_{mir} = \alpha_m + \mu_i + c_{mi} + \epsilon_{mir}, \quad e_{mi} \sim \mathcal{N}(0, \sigma_m^2), \quad c_{mi} \sim \mathcal{N}(0, \tau_m^2). \quad (4.12)$$

Let y_{mir} be the r th replicate measurement on the i th item by the m th method, where $m = 1, 2$, $i = 1, \dots, N$, and $r = 1, \dots, n_i$. When the design is balanced and there is no ambiguity we can set

$n_i = n$. The LME model can be written

$$y_{mir} = \beta_0 + \beta_m + b_{mi} + \epsilon_{mir}. \quad (4.13)$$

Here β_0 and β_m are fixed-effect terms representing, respectively, a model intercept and an overall effect for method m . The model can be reparameterized by gathering the β terms together into (fixed effect) intercept terms $\alpha_m = \beta_0 + \beta_m$. The b_{1i} and b_{2i} terms are correlated random effect parameters having $E(b_{mi}) = 0$ with $\text{Var}(b_{mi}) = d_m^2$ and $\text{Cov}(b_{1i}, b_{2i}) = d_{12}$.

The random error term for each response is denoted ϵ_{mir} having $E(\epsilon_{mir}) = 0$, $\text{Var}(\epsilon_{mir}) = \sigma_m^2$, $\text{Cov}(b_{mir}, b_{m'ir}) = \sigma_{12}$, $\text{Cov}(\epsilon_{mir}, \epsilon_{mir'}) = 0$ and $\text{Cov}(\epsilon_{mir}, \epsilon_{m'ir'}) = 0$. When two methods of measurement are in agreement, there is no significant differences between β_1 and β_2 , d_1^2 and d_2^2 , and σ_1^2 and σ_2^2 .

Additionally these parameter are assumed to have Gaussian distribution. Two methods of measurement are in complete agreement if the null hypotheses $H_1: \alpha_1 = \alpha_2$ and $H_2: \sigma_1^2 = \sigma_2^2$ and $H_3: d_1^2 = d_2^2$ hold simultaneously. Roy (2009) uses a Bonferroni correction to control the familywise error rate for tests of $\{H_1, H_2, H_3\}$ and account for difficulties arising due to multiple testing. Additionally, Roy combines H_2 and H_3 into a single testable hypothesis $H_4: \omega_1^2 = \omega_2^2$, where $\omega_m^2 = \sigma_m^2 + d_m^2$ represent the overall variability of method m .

Here the terms α_m and μ_i represent the fixed effect for method m and a true value for item i respectively. The random effect terms comprise an interaction term c_{mi} and the residuals ϵ_{mir} . The c_{mi} term represent random effect parameters corresponding to the two methods, having $E(c_{mi}) = 0$ with $\text{Var}(c_{mi}) = \tau_m^2$.

Carstensen specifies the variance of the interaction terms as being univariate normally distributed. As such, $\text{Cov}(c_{mi}, c_{m'i}) = 0$. All the random effects are assumed independent, and that all replicate measurements are assumed to be exchangeable within each method.

The presence of the true value term μ_i gives rise to an important difference between Carstensen's and Roy's models. Of particular importance is terms of the model, a true value for item i (μ_i). The fixed effect of Roy's model comprise of an intercept term and fixed effect terms for both methods, with no reference to the true value of any individual item. In other words, Roy considers the group of items being measured as a sample taken from a population. Therefore a distinction can be made between the

two models: Roy's model is a standard LME model, whereas Carstensen's model is a more complex additive model.

Chapter 5

Bradley Blackwood

5.1 Bartko's Bradley-Blackwood Test

This is a regression based approach that performs a simultaneous test for the equivalence of means and variances of the respective methods. We have identified this approach to be examined to see if it can be used as a foundation for a test perform a test on means and variances individually.

$$D = (X_1 - X_2) \quad (5.1)$$

$$M = (X_1 + X_2)/2 \quad (5.2)$$

The Bradley Blackwood Procedure fits D on M as follows:

$$D = \beta_0 + \beta_1 M \quad (5.3)$$

- The Bradley Blackwood test is a simultaneous test for bias and precision. They propose a regression approach which fits D on M, where D is the difference and average of a pair of results.
- Both beta values, the intercept and slope, are derived from the respective means and standard deviations of their respective data sets.
- We determine if the respective means and variances are equal if both beta values are simultaneously equal to zero. The Test is conducted using an F test, calculated from the results of a regression of D on M.

- We have identified this approach to be examined to see if it can be used as a foundation for a test perform a test on means and variances individually.
- Russell et al have suggested this method be used in conjunction with a paired t-test , with estimates of slope and intercept.

5.2 Bradley-Blackwood Test (Kevin Hayes Talk)

This work considers the problem of testing $\mu_1 = \mu_2$ and $\sigma_1^2 = \sigma_2^2$ using a random sample from a bivariate normal distribution with parameters $(\mu_1, \mu_2, \sigma_1^2, \sigma_2^2, \rho)$.

The new contribution is a decomposition of the Bradley-Blackwood test statistic (*Bradley and Blackwood, 1989*) for the simultaneous test of $\mu_1 = \mu_2$; $\sigma_1^2 = \sigma_2^2$ as a sum of two statistics.

One is equivalent to the Pitman-Morgan (*Pitman, 1939; Morgan, 1939*) test statistic for $\sigma_1^2 = \sigma_2^2$ and the other one is a new alternative to the standard paired-t test of $\mu_D = \mu_1 - \mu_2 = 0$.

Surprisingly, the classic Student paired-t test makes no assumptions about the equality (or otherwise) of the variance parameters.

The power functions for these tests are quite easy to derive, and show that when $\sigma_1^2 = \sigma_2^2$, the paired t-test has a slight advantage over the new alternative in terms of power, but when $\sigma_1^2 \neq \sigma_2^2$, the new test has substantially higher power than the paired-t test.

While Bradley and Blackwood provide a test on the joint hypothesis of equal means and equal variances their regression based approach does not separate these two issues.

The rejection of the joint hypothesis may be due to two groups with unequal means and unequal variances; unequal means and equal variances, or equal means and unequal variances. We propose an approach for resolving this (model selection) problem in a manner controlling the magnitudes of the relevant type I error probabilities.

Deming Regression

- Informative analysis for the purposes of method comparison, Deming Regression is a regression technique taking into account uncertainty in both the independent and dependent variables.

- Demings method always results in one regression fit, regardless of which variable takes the place of the predictor variables.
- The measurement error (λ or λ) is specified with measurement error variance related as

$$\lambda = \sigma_y^2 / \sigma_x^2$$

(where σ_x^2 and σ_y^2 is the measurement error variance of the x and y variables, respectively).

- In the case where λ is equal to one, (i.e. equal error variances), the methodology is equivalent to *orthogonal regression*.
- Deming approaches the matter by simultaneously minimizing the sum of the square of the residuals of both variables. This derivation results in the best fit to minimize the sum of the squares of the perpendicular distances from the data points.
- To compute the slope by Demings formula, normally distributed error of both variables is assumed, as well as a constant level of imprecision throughout the range of measurements.

5.3 Simple Linear Regression

Simple linear regression is defined as such with the name ‘Model I regression’ by Cornbleet Gochman (1979), in contrast to ‘Model II regression’.

On account of the fact that one set of measurements are linearly related to another, one could surmise that Linear Regression is the most suitable approach to analyzing comparisons. This approach is unsuitable on two counts. Firstly one of the assumptions of Regression analysis is that the independent variable values are without error. In method comparison studies one must assume the opposite; that there is error present in the measurements. Secondly a regression of X on Y would yield an entirely different result from Y on X .

Simple linear regression calculates a line of best fit for two sets of data, in which the independent variable, X , is measured without error, with y as the dependent variable.

SLR (Model I) regression is considered by many Altman and Bland (1983); Cornbleet and Cochrane (1979); Ludbrook (1997) to be wholly unsuitable for method comparison studies, although recommended for use in calibration studies [Corncoch]. Even in the case where one method is a gold standard, it is disputed as to whether it is a valid approach. Model II regression is more suitable for method comparison studies, but it is more difficult to execute. Both Model I and II regression models are unduly influenced by outliers. Regression Models can not be used to analyze repeated measurements

Regression Analysis

Another inappropriate approach is the regressing one set of measurements against the other. According to this methodology the measurement methods could be considered equivalent if the confidence interval for the regression coefficient included 1. Analysts sometimes use least squares (referred to by Ludbrook as Model I) regression analysis to calibrate one method of measurement against another. In this technique, the sum of the squares of the vertical deviations of y values from the line is minimized. This approach is invalid, because both y and x values are attended by random error.

The Identity Plot

This is a simple graphical approach, advocated by Bland and Altman (1986), that yields a cursory examination of how well the measurement methods agree. In the case of good agreement, the co-variates of the plot accord closely with the $X = Y$ line.

Advantages of Regression Approaches for MCS

- These methods can be employed in conversion problems.
- Bland and Altman have stated that regression analysis offers insights into MCS problems.

Disadvantages

- Regression methods are uninformative about the variability of the differences.
- Regression methods can determine the presence of bias, and the levels of constant bias and proportional bias thereof Ludbrook (1997, 2002).

5.4 Constant and Proportional Bias

Linear Regression is a commonly used technique for comparing paired assays. The Intercept and Slope can provide estimates for the constant bias and proportional bias occurring between both methods. If the basic assumptions underlying linear regression are not met, the regression equation, and consequently the estimations of bias are undermined. Outliers are a source of error in regression estimates.

Constant or proportional bias in method comparison studies using linear regression can be detected by an individual test on the intercept or the slope of the line regressed from the results of the two methods to be compared.

Bartko's Discussion of BB

Let $y = X_1 - X_2$ and $x = (X_1 + X_2)/2$. The Bradley-Blackwood procedure fits y on x , such that

$$y = \beta_0 + \beta_1 x$$

The slope and intercepts are given by

$$\beta_1 = \frac{(\sigma_1^2 - \sigma_2^2)}{2\sigma_x^2}$$

Pitman's Test on Correlated variances

$$H_0 : \sigma_1^2 = \sigma_2^2$$

$$H_0 : \sigma_1^2 = \sigma_2^2$$

Pitman's test is identical to the slope equal to zero in the regression of y on x .

5.5 Conclusions about Existing Methodologies

The Bland Altman methodology is well noted for its ease of use, and can be easily implemented with most software packages. Also it doesn't require the practitioner to have more than basic statistical training.

The plot is quite informative about the variability of the differences over the range of measurements. For example, an inspection of the plot will indicate the 'fan effect'. They also can be used to detect the presence of an outlier.

Ludbrook (1997, 2002) criticizes these plots on the basis that they presents no information on effect of constant bias or proportional bias. These plots are only practicable when both methods measure in the same units. Hence they are totally unsuitable for conversion problems. The limits of agreement are somewhat arbitrarily constructed. They may or may not be suitable for the data in question. It has been found that the limits given are too wide to be acceptable. There is no guidance on how to deal with outliers. Bland and Altman recognize effect they would have on the limits of agreement, but offer no guidance on how to correct for those effects.

There is no formal testing procedure provided. Rather, it is upon the practitioner opinion to judge the outcome of the methodology.

Bartko's Ellipse

$$\frac{(x - \bar{x})^2}{\sigma_x^2} - \frac{2\rho(x - \bar{x})(y - \bar{y})}{\sigma_x\sigma_y} + \frac{(y - \bar{y})^2}{\sigma_y^2} = \chi^2(2df(1 - \rho^2))$$

section*Remarks

- Pearson's Correlation of (x,y) is the same as Pitman's correlation of sums and differences.
- Techniques for plotting an ellipse can be found in Douglas Altman's book.

5.6 A regression based approach based on Bland Altman Analysis

Bland and Altman have stated that regression analysis offers insights into method comparison studies. Regression methods can determine the presence of bias, and the levels of constant bias and proportional bias thereof Ludbrook (1997, 2002). While they are informative about inter-method bias, Regression methods offer the analyst no insights into the relative precision of both methods. These methods can be employed in conversion problems, however errors are attended. *Lu et al* used such a technique in their comparison of DXA scanners. They also used the Blackwood Bradley test. However it was shown that, for particular comparisons, agreement between methods was indicated according to one test, but lack of agreement was indicated by the other.

Remarks

- Pearson's Correlation of (x,y) is the same as Pitman's correlation of sums and differences.
- Techniques for plotting an ellipse can be found in Douglas Altman's book.

5.7 The MCR R package - Regression Techniques for MCS

The *mcr* packages provides a set of regression techniques to quantify the relation between two measurement methods.

In particular, it address regression problems with errors in both variables, but without repeated measurements. The *mcr* package follows the CLSI EP09-A3 recommendations for analytical method comparison and estimation of bias using patient samples.

Methods featured in the mcr package

- Deming Regression
- Weighted Deming Regression

- Passing-Bablok Regression

The *creatinine* gives the blood and serum preoperative creatinine measurements in 110 heart surgery patients.

```
library("mcr")
data("creatinine", package="mcr")
tail(creatinine)

fit.lr <- mcreg(as.matrix(creatinine), method.reg="LinReg", na.rm=TRUE)
fit.wlr <- mcreg(as.matrix(creatinine), method.reg="WLinReg", na.rm=TRUE)
compareFit( fit.lr, fit.wlr )
```

5.8 Implementation of Deming Regression with Rs

Thus far, one of the few R implementations of Deming regression is contained in the ‘MethComp’ package. (Carstensen et al., 2008).

Unless specified otherwise, the variance ratio λ has a default value of one. A means of computing likelihood functions would potentially allow for an algorithm for estimating the true variance ratio.

5.9 KP

Most residual covariance structures are design for one within-subject factor. However two or more may be present. For such cases, an appropriate approach would be the residual covariance structure using Kronecker product of the underlying within-subject factor specific covariances structure.

Chapter 6

Residual Diagnostics

The original Bland Altman Method was developed for two sets of measurements done on one occasion (i.e. independent data), and so this approach is not suitable for repeated measures data. However, as a naive analysis, it may be used to explore the data because of the simplicity of the method. Myles states that such misuse of the standard Bland Altman method is widespread in Anaesthetic and critical care literature.

Bland and Altman have provided a modification for analysing repeated measures under stable or changing conditions, where repeated data is collected over a period of time. Myers proposes an alternative Random effects model for this purpose.

with repeated measures data, we can calculate the mean of the repeated measurements by each method on each individual. *The pairs of means can then be used to compare the two methods based on the 95% limits of agreement for the difference of means. The bias between the two methods will not be affected by averaging the repeated measurements.* However the variation of the differences will be underestimated by this practice because the measurement error is, to some extent, removed. Some advanced statistical calculations are needed to take into account these measurement errors. *Random effects models can be used to estimate the within-subject variation after accounting for other observed and unobserved variations, in which each subject has a different intercept and slope over the observation period. On the basis of the within-subject variance estimated by the random effects model, we can then create an appropriate*

Bland Altman Plot. The sequence or the time of the measurement over the observation period can be taken as a random effect.

Taxonomy of Cook's Distances for LMEs

Zewotir and Galpin (2005) discusses a taxonomy of Cook's distance when applied to LME models.

- For variance components γ : $CD(\gamma)_i$,
- For fixed effect parameters β : $CD(\beta)_i$,
- For random effect parameters \mathbf{u} : $CD(u)_i$,
- For linear functions of $\hat{\beta}$: $CD(\psi)_i$

Computational Limitations for Cook's Distance

Application of Cook's Distances are limited by computation tractability.

Application of case-deletion diagnostics offer some interested for Method Comparison Studies

Care must be given when interpreting these plots. For example the position of case 68 on the BSVR indicates that that case 68

Any diagnostic plot may constructed using Overall variability and intermethod bias.

6.0.1 Case-Deletion Diagnostics

Influence diagnostics are formal techniques that assess the influence of observations on parameter estimates for β and θ . A common technique is to refit the model with an observation or group of observations omitted.

Case-deletion diagnostics provide a useful tool for identifying influential observations and outliers. Since the pioneering work of Cook in 1977, deletion measures have been applied to many statistical models for identifying influential observations.

Standard statistical packages concentrate on calculating and testing parameter estimates without considering the diagnostics of the model.

The assessment of the effects of perturbations in data, on the outcome of the analysis, is known as statistical influence analysis. Influence analysis examines the robustness of the model.

6.1 Cook's Distance

As well as individual observations, Cook's distance can be used to analyse the influence of observations in subset U on a vector of parameter estimates (Cook, 1977).

$$\hat{e}_{i(U)} = y_i - x\hat{\beta}_{(U)} \quad (6.1)$$

$$\delta_{(U)} = \hat{\beta} - \hat{\beta}_{(U)} \quad (6.2)$$

6.2 Haslett Dillane Hayes

Haslett and Dillane (2004) offers an procedure to assess the influences for the variance components within the linear model, complementing the existing methods for the fixed components. The essential problem is that there is no useful updating procedures for \hat{V} , or for \hat{V}^{-1} .

Haslett and Dillane (2004) remark that linear mixed effects models didn't experience a corresponding growth in the use of deletion diagnostics, adding that McCullough and Searle (2001) makes no mention of diagnostics whatsoever.

Haslett and Dillane (2004) propose an alternative, and computationally inexpensive approach, making use of the 'delete=replace' identity.

Haslett (1999) considers the effect of 'leave k out' calculations on the parameters β and σ^2 , using several key results from Haslett and Hayes (1998) on partitioned matrices.

6.3 Demidenk Case Deletion Diagnostics

Demidenko (2004) extends several regression diagnostic techniques commonly used in linear regression, such as leverage, infinitesimal influence, case deletion diagnostics, Cook's distance, and local influence to the linear mixed-effects model. In each case, the proposed new measure has a direct interpretation

in terms of the effects on a parameter of interest, and reduces to the familiar linear regression measure when there are no random effects.

The new measures that are proposed by Demidenko (2004) are explicitly defined functions and do not require re-estimation of the model, especially for cluster deletion diagnostics. The basis for both the cluster deletion diagnostics and Cook's distance is a generalization of Miller's simple update formula for case deletion for linear models. Furthermore Demidenko (2004) shows how Pregibon's infinitesimal case deletion diagnostics is adapted to the linear mixed-effects model.

6.4 Cook's Distance - Implementation with R

Cook's Distance is a measure indicating to what extent model parameters are influenced by (a set of) influential data on which the model is based. This function computes the Cook's distance based on the information returned by the `estex()` function.

6.5 LME diagnostic measures

6.5.1 Andrews-Pregibon statistic

- For fixed effect parameters β .

The Andrews-Pregibon statistic AP_i is a measure of influence based on the volume of the confidence ellipsoid. The larger this statistic is for observation i , the stronger the influence that observation will have on the model fit.

6.5.2 Cook's Distance

- For variance components γ

Diagnostic tool for variance components

$$C_{\theta i} = ((\hat{\theta})_{[i]} - \hat{\theta})^T \text{cov}(\hat{\theta})^{-1} ((\hat{\theta})_{[i]} - \hat{\theta})$$

6.5.3 Variance Ratio

- For fixed effect parameters β .

6.5.4 Cook-Weisberg statistic

- For fixed effect parameters β .

6.5.5 Andrews-Pregibon statistic

- For fixed effect parameters β .

The Andrews-Pregibon statistic AP_i is a measure of influence based on the volume of the confidence ellipsoid. The larger this statistic is for observation i , the stronger the influence that observation will have on the model fit.

6.6 Two-tailed testing

A test for equality of variances, based on the likelihood Ratio test, is very simple to implement using existing methodologies. All that is required it to specify the reference model and the relevant nested mode as arguments to the command `anova()`. The output can be interpreted in the usual way.

6.7 One Tailed Testing

The approach proposed by Roy deals with the question of agreement, and indeed interchangeability, as developed by Bland and Altman's corpus of work. In the view of Dunn, a question relevant to many practitioners is which of the two methods is more precise.

The relationship between precision and the within-item and between-item variability must be established. Roy establishes the equivalence of repeatability and within-item variability, and hence precision. The method with the smaller within-item variability can be deemed to be the more precise.

6.8 Enabling One Tailed Testing

A useful approach is to compute the confidence intervals for the ratio of within-item standard deviations (equivalent to the ratio of repeatability coefficients), which can be interpreted in the usual manner (or alternatively, the ratio of the variances). In fact, the ratio of within-item standard deviations, with the attendant confidence interval, can be determined using a single R command: `intervals()`.

Pinheiro and Bates (pg 93-95) give a description of how confidence intervals for the variance components are computed. Furthermore a complete set of confidence intervals can be computed to complement the variance component estimates. However , to facilitate one tailed testing, What is required is the computation of the variance ratios of within-item and between-item standard deviations.

A nave approach would be to compute the variance ratios by relevant F distribution quantiles. However, the question arises as to the appropriate degrees of freedom. However, Douglas Bates has stated that an alternative approach is required (i.e. Profile Likelihoods)

”The omission of standard errors on variance components is intentional. The distribution of an estimator of a variance component is highly skewed and obtaining an estimate of the standard deviation of a skewed distribution is not very useful. A much better approach is based on profiling the objective function.” (Douglas Bates May 2012)

6.9 Profile Likelihood

Normal-based confidence intervals for a parameter of interest are inaccurate when the sampling distribution of the estimate is skewed. The technique known as profile likelihood can produce confidence intervals with better coverage. It may be used when the model includes only the variable of interest or several other variables in addition. Profile-likelihood confidence intervals are particularly useful in nonlinear models.

Profile likelihood confidence intervals are based on the log-likelihood function.

6.10 Implementation of PL Confidence Intervals

The suitable calculation of confidence limits for this variance ratio are to be computed using the profile likelihood approach. The R package `profilelikelihood` will be assessed for feasibility, particularly the command `profilelikelihood.lme()`

Normal-based con

dence intervals for a parameter of interest are inaccurate when the sampling distribution of the estimate is skewed. The technique known as profile likelihood can produce confidence intervals with better coverage. It may be used when the model includes only the variable of interest or several other variables in addition. Profile-likelihood confidence intervals are particularly useful in nonlinear models. Profile likelihood confidence intervals are based on the log-likelihood function.

6.11 Zewotir: Computation and Notation

with \mathbf{V} unknown, a standard practice for estimating $\mathbf{X}\boldsymbol{\beta}$ is the estimate the variance components σ_j^2 , compute an estimate for \mathbf{V} and then compute the projector matrix A , $\mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{A}\mathbf{Y}$.

Zewotir remarks that \mathbf{D} is a block diagonal with the i -th block being $u\mathbf{I}$

6.12 Haslett Hayes

For fixed effect linear models with correlated error structure Haslett (1999) showed that the effects on the fixed effects estimate of deleting each observation in turn could be cheaply computed from the fixed effects model predicted residuals.

A general theory is presented for residuals from the general linear model with correlated errors. It is demonstrated that there are two fundamental types of residual associated with this model, referred to here as the marginal and the conditional residual. These measure respectively the distance to the global aspects of the model as represented by the expected value and the local aspects as represented by the conditional expected value. These residuals may be multivariate.

In contrast to classical linear models, diagnostics for LME are difficult to perform and interpret,

because of the increased complexity of the model

6.13 Confounded Residuals

Hilden-Minton (1995, PhD thesis, UCLA): residual is pure for a specific type of error if it depends only on the fixed components and on the error that it is supposed to predict. Residuals that depend on other types of errors are called ***confounded residuals***. This code will allow you to make QQ plots for each level of the random effects. LME models assume that not only the within-cluster residuals are normally distributed, but that each level of the random effects are as well. Depending on the model, you can vary the level from 0, 1, 2 and so on.

```
qqnorm(JS.roy1, ~ranef(.))  
  
# qqnorm(JS.roy1, ~ranef(.,levels=1))
```

This code will allow you to make QQ plots for each level of the random effects. LME models assume that not only the within-cluster residuals are normally distributed, but that each level of the random effects are as well. Depending on the model, you can vary the level from 0, 1, 2 and so on.

```
qqnorm(JS.roy1, ~ranef(.))  
  
# qqnorm(JS.roy1, ~ranef(.,levels=1))
```

Chapter 7

Fitting LME Models

Further to previous material, an appraisal of the current state of development for statistical software for fitting for LME models, particularly for `nlme` and `lme4` fitted models.

The **`lme4`** package is used to fit linear and generalized linear mixed-effects models in the R environment. The **`lme4`** package is also under active development, under the leadership of Ben Bolker (McMaster Uni., Canada).

Crucially, a review of internet resources indicates that almost all of the progress in this regard has been done for `lme4` fitted models, specifically the *Influence.ME* R package. (Nieuwenhuis et al 2014) Conversely there is very little for `nlme` models. One would immediately look at the current development workflow for both packages.

As an aside, Douglas Bates was arguably the most prominent R developer working in the LME area. However Bates has now prioritised the development of LME models in another computing environment, i.e Julia.

With regards to `nlme`, the package is now maintained by the R core development team. The most recent major text is by Galecki & Burzykowski, who have published *Linear Mixed Effects Models using R*. Also, the accompanying R package, `nlmeU` package is under current development, with a version being released 0.70 – 3.

7.1 Relevance of Roy's Methodology

The relevance of Roy's methodology is that estimates for the between-item variances for both methods \hat{d}_m^2 are computed. Also the VC matrices are constructed with covariance terms and, so the difference variance must be formulated accordingly.

$$\hat{\alpha}_1 - \hat{\alpha}_2 \pm \sqrt{\hat{d}_1^2 + \hat{d}_2^2 + \hat{\sigma}_1^2 + \hat{\sigma}_2^2 - 2\hat{d}_{12} - 2\hat{\sigma}_{12}}$$

7.2 Interaction Terms in Model

One important feature of replicate observations is that they should be independent of each other. In essence, this is achieved by ensuring that the observer makes each measurement independent of knowledge of the previous value(s). This may be difficult to achieve in practice.

Further to Barnhart et al. (2007), if the measurements by a method on an item are not necessarily true replications, e.g., repeated measures over time, then additional terms may be needed for e_{mir} . Carstensen et al. (2008) also addresses this issue by the addition of an interaction term (i.e. a random effect) u_{mi} , yielding

$$y_{mir} = \alpha_{mi} + u_{mi} + e_{mi}.$$

The additional interaction term is characterized as $u_{mi} \sim \mathcal{N}(0, \tau_m^2)$ (Carstensen et al., 2008). This extra interaction term provides a source of extra variability, but this variance is not relevant to computing the case-wise differences.

7.3 Difference Variance further to Carstensen

Even though the separate variances can not be identified, their sum can be estimated by the empirical variance of the differences.

Like wise the separate α can not be estimated, only their difference can be estimated as \bar{D}

We assume that the variance of the measurements is different for both methods, but it does not mean that the separate variances can be estimated with the data available.

7.4 Why use LMEs for Method Comparison?

The LME model approach has seen increased use as a framework for method comparison studies in recent years (Lai & Shaio, Carstensen and Choudhary as examples). In part this is due to the increased profile of LME models, and furthermore the availability of capable software. Additionally LME based approaches may utilise the diagnostic and influence analysis techniques that have been developed in recent times.

Roy proposes an LME model with Kronecker product covariance structure in a doubly multivariate setup. Response for i th subject can be written as

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + b_{1i} z_{i1} + b_{2i} z_{i2} + \epsilon_i$$

- β_1 and β_2 are fixed effects corresponding to both methods. (β_0 is the intercept.)
- b_{1i} and b_{2i} are random effects corresponding to both methods.

Overall variability between the two methods (Ω) is sum of between-subject (D) and within-subject variability (Σ),

$$\text{Block } \Omega_i = \begin{bmatrix} d_1^2 & d_{12} \\ d_{12} & d_2^2 \end{bmatrix} + \begin{bmatrix} \sigma_1^2 & \sigma_{12} \\ \sigma_{12} & \sigma_2^2 \end{bmatrix}.$$

The well-known “Limits of Agreement”, as developed by Bland and Altman (1986) are easily computable using the LME framework, proposed by Roy. While we will not be considering this analysis, a demonstration will be provided in the example.

Further to this, Roy(2009) demonstrates an suite of tests that can be used to determine how well two methods of measurement, in the presence of repeated measures, agree with each other.

- No Significant inter-method bias
- No difference in the between-subject variabilities of the two methods

- No difference in the within-subject variabilities of the two methods

7.5 Definition of Replicate measurements

Further to Bland and Altman (1999), a formal definition is required of what exactly replicate measurements are

By replicates we mean two or more measurements on the same individual taken in identical conditions. In general this requirement means that the measurements are taken in quick succession.

Bland and Altman (1999) also remark that an important feature of replicate observations is that they should be independent of each other. This issue is addressed by Carstensen (2010), in terms of exchangeability and linkage. Carstensen advises that repeated measurements come in two *substantially different* forms, depending on the circumstances of their measurement: exchangeable and linked.

7.5.1 Exchangeable measurements

Repeated measurements are said to be exchangeable if no relationship exists between successive measurements across measurements. If the condition of exchangeability exists, a group of measurement of the same item determined by the same method can be re-arranged in any permutation without prejudice to proper analysis. There is no reason to believe that the true value of the underlying variable has changed over the course of the measurements.

For the purposes of method comparison studies the following remarks can be made. The r -th measurement made by method 1 has no special correspondence to the r -th measurement made by method 2, and consequently any pairing of repeated measurements are as good as each other.

Exchangeable repeated measurements can be treated as true replicates.

7.5.2 Linked measurements

Repeated measurements are said to be linked if a direct correspondence exists between successive measurements across measurements, i.e. pairing. Such measurements are commonly made with a time interval between them, but simultaneously for both methods. Paired measurements are exchangeable,

but individual measurements are not.

If the paired measurements are taken in a short period of time so that no real systemic changes can take place on each item, they can be considered true replicates. Should enough time elapse for systemic changes, linked repeated measurements can not be treated as true replicates.

7.5.3 Replicate measurements in ARoy2009's paper

Roy (2009) takes its definition of replicate measurement: two or more measurements on the same item taken under identical conditions. ARoy2009 also assumes linked measurements, but it is can be used for the non-linked case.

7.5.4 Random effects

Further to Barnhart et al. (2007), if the measurements by a method on an item are not necessarily true replications, e.g., repeated measures over time, then additional terms may be needed for e_{mir} . Carstensen et al. (2008) also addresses this issue by the addition of an interaction term (i.e. a random effect) u_{mi} , yielding

$$y_{mir} = \alpha_{mi} + u_{mi} + e_{mi}.$$

The additional interaction term is characterized as $u_{mi} \sim \mathcal{N}(0, \tau_m^2)$ (Carstensen et al., 2008).

This extra interaction term provides a source of extra variability, but this variance is not relevant to computing the case-wise differences.

Carstensen et al. (2008) advises that the formulation of the model should take the exchangeability (in other words, whether or not the measurements are ‘true replicates’) into account. If there is a linkage between measurements (therefore not ‘true’ replicates), the ‘item by replicate’ should be included in the model. If there is no linkage, and the replicates are indeed true replicates, the interaction term should be omitted.

Carstensen et al. (2008) demonstrates how to compute the limits of agreement for two methods in the case of linked measurements. As a surplus source of variability is excluded from the computation, the limits of agreement are not unduly wide, which would have been the case if the measurements were treated as true replicates.

Roy (2009) also assigns a random effect u_{mi} for each response y_{mir} . Importantly ARoy2009’s model assumes linkage.

7.6 Model for replicate measurements

We generalize the single measurement model for the replicate measurement case, by additionally specifying replicate values. Let y_{mir} be the r –th replicate measurement for subject “i” made by method “m”. Further to Barnhart et al. (2007) fixed effect can be expressed with a single term α_{mi} , which incorporate the true value μ_i .

$$y_{mir} = \mu_i + \alpha_m + e_{mir}$$

Combining fixed effects (Barnhart et al., 2007), we write,

$$y_{mir} = \alpha_{mi} + e_{mir}.$$

The following assumptions are required

- e_{mir} is independent of the fixed effects with mean $E(e_{mir}) = 0$.
- Further to Barnhart et al. (2007) between-item and within-item variances $\text{Var}(\alpha_{mi}) = \sigma_{Bm}^2$ and $\text{Var}(e_{mir}) = \sigma_{Wm}^2$
- In keeping with Roy (2009), these variance shall be considered as part of the between-item variance covariance matrix \mathbf{D} and the within-item variance covariance matrix $\mathbf{\Sigma}$ respectively, and will be denoted accordingly (i.e. d_m^2 and σ_m^2).
- Additionally, the total variability of method "m", denoted ω_m^2 is the sum of the within-item and between-item variabilities.

$$\omega_m^2 = d_m^2 + \sigma_m^2$$

Chapter 8

BA99

8.1 Regression-based Limits of Agreement

Assuming that there will be no curvature in the scatter-plot, the methodology regresses the difference of methods (d) on the average of those methods (a) with a simple intercept slope model; $\hat{d} = b_0 + b_1 a$. Should the slope b_1 be found to be negligible, \hat{d} takes the value \bar{d} .

The next step to take in calculating the limits is also a regression, this time of the residuals as a function of the scale of the measurements, expressed by the averages a_i ; $\hat{R} = c_0 + c_1 a_i$

With reference to absolute values following a half-normal distribution with mean $\sigma\sqrt{\frac{2}{\pi}}$, Bland and Altman (1999) formulate the regression based limits of agreement as follows

$$\hat{d} \pm 1.96\sqrt{\frac{\pi}{2}}\hat{R} = \hat{d} \pm 2.46\hat{R} \quad (8.1)$$

8.2 Steps of Structural Equation modelling

1. **Model Specification** We must state the theoretical model either as a set of equations.
2. **Identification** This step involves checking that the model can be estimated with observable data, both in theory and in practice.
3. **Estimation** The models parameters are statistically estimated from data. (multiple regression is one such method)
4. **Model Fit** The estimated model parameters are used to predict the correlations and covariance between measured variables The predicted correlations, or covariance are compared to the observed correlations, or covariance. (Measures of model fit are calculated)

Chapter 9

Appendices 1

Chapter 10

Augmented GLMs

Generalized linear models are a generalization of classical linear models.

10.1 Augmented GLMs

With the use of h-likelihood, a random effected model of the form can be viewed as an ‘augmented GLM’ with the response variables $(y^t, \phi_m^t)^t$, (with $\mu = E(y), u = E(\phi), \text{var}(y) = \theta V(\mu)$). The augmented linear predictor is

$$\eta_{ma} = (\eta^t, \eta_m^t)^t = T\omega.$$

.

The subscript M is a label referring to the mean model.

$$\begin{pmatrix} Y \\ \psi_M \end{pmatrix} = \begin{pmatrix} X & Z \\ 0 & I \end{pmatrix} \begin{pmatrix} \beta \\ \nu \end{pmatrix} + e^* \quad (10.1)$$

The error term e^* is normal with mean zero. The variance matrix of the error term is given by

$$\Sigma_a = \begin{pmatrix} \Sigma & 0 \\ 0 & D \end{pmatrix}. \quad (10.2)$$

$$y_a = T\delta + e^*$$

Weighted least squares equation

10.1.1 The Augmented Model Matrix

$$X = \begin{pmatrix} T & Z \\ 0 & I \end{pmatrix} \delta = \begin{pmatrix} \beta \\ \nu \end{pmatrix} \quad (10.3)$$

newpage

10.2 Algorithms : ML v REML

Maximum likelihood estimation is a method of obtaining estimates of unknown parameters by optimizing a likelihood function. The ML parameter estimates are the values of the argument that maximise the likelihood function, i.e. the estimates that make the observed values of the dependent variable most likely, given the distributional assumptions

The most common iterative algorithms used for the optimization problem in the context of LMEs are the EM algorithm, fisher scoring algorithm and NR algorithm, which [cite:West] commends as the preferred method.

A mixed model is an extension of the general linear models that can specify additional random effects terms.

Parameter of the mixed model can be estimated using either ML or REML, while the AIC and the BIC can be used as measures of "goodness of fit" for particular models, where smaller values are considered preferable.

(*Wikipedia*)The restricted (or residual, or reduced) maximum likelihood (REML) approach is a particular form of maximum likelihood estimation which does not base estimates on a maximum likelihood fit of all the information, but instead uses a likelihood function calculated from a transformed set of data, so that nuisance parameters have no effect.

In contrast to the earlier maximum likelihood estimation, REML can produce unbiased estimates of variance and covariance parameters.

ML procedures for LME

The maximum likelihood procedure of Hartley and Rao yields simultaneous estimates for both the fixed effects and the random effect, by maximising the likelihood of \mathbf{y} with respect to each element of $\boldsymbol{\beta}$ and \mathbf{b} .

10.3 Estimation of random effects

Estimation of random effects for LME models in the NLME package is accomplished through use of both EM (Expectation-Maximization) algorithms and Newton-Raphson algorithms.

- EM iterations bring estimates of the parameters into the region of the optimum very quickly, but convergence to the optimum is slow when near the optimum.
- Newton-Raphson iterations are computationally intensive and can be unstable when far from the optimum. However, close to the optimum they converge quickly.
- The LME function implements a hybrid approach, using 25 EM iterations to quickly get near the optimum, then switching to Newton-Raphson iterations to quickly converge to the optimum.
- If convergence problems occur, the “controlargument in LME can be used to change the way the model arrives at the optimum.

10.4 Covariance Parameters

The unknown variance elements are referred to as the covariance parameters and collected in the vector θ .

10.4.1 Methods and Measures

The key to making deletion diagnostics useable is the development of efficient computational formulas, allowing one to obtain the case deletion diagnostics by making use of basic building blocks, computed only once for the full model.

Zewotir and Galpin (2005) lists several established methods of analyzing influence in LME models. These methods include

- Cook’s distance for LME models,
- likelihood distance,

- the variance (information) ration,
- the Cook-Weisberg statistic,
- the Andrews-Prebigon statistic.

10.5 Haslett's Analysis

For fixed effect linear models with correlated error structure Haslett (1999) showed that the effects on the fixed effects estimate of deleting each observation in turn could be cheaply computed from the fixed effects model predicted residuals.

10.6 Computation and Notation

with \mathbf{V} unknown, a standard practice for estimating $\mathbf{X}\boldsymbol{\beta}$ is to estimate the variance components σ_j^2 , compute an estimate for \mathbf{V} and then compute the projector matrix A , $\mathbf{X}\hat{\boldsymbol{\beta}} = A\mathbf{Y}$.

Zewotir and Galpin (2005) remarks that \mathbf{D} is a block diagonal with the i -th block being $u\mathbf{I}$

Chapter 11

Generalized linear models

11.1 Generalized Linear model

In statistics, the generalized linear model (GzLM) is a flexible generalization of ordinary least squares regression. The GzLM generalizes linear regression by allowing the linear model to be related to the response variable via a link function and by allowing the magnitude of the variance of each measurement to be a function of its predicted value.

Mixed Effects Models offer a flexible framework by which to model the sources of variation and correlation that arise from grouped data. This grouping can arise when data collection is undertaken in a hierarchical manner, when a number of observations are taken on the same observational unit over time, or when observational units are in some other way related, violating assumptions of independence.

11.2 Generalized Model(GzLM)

Nelder and Wedderburn (1972) integrated the previously disparate and separate approaches to models for non-normal cases in a framework called "generalized linear models." The key elements of their approach is to describe any given model in terms of its link function and its variance function.

11.2.1 What is a GzLM

$$E(\mathbf{Y}) = \boldsymbol{\mu} = g^{-1}(\mathbf{X}\boldsymbol{\beta}) \quad (11.1)$$

where $E(Y)$ is the expected value of Y , $X\beta$ is the linear predictor, a linear combination of unknown parameters, β and g is the link function.

$$\text{Var}(\mathbf{Y}) = V(\boldsymbol{\mu}) = V(g^{-1}(\mathbf{X}\boldsymbol{\beta}))$$

11.2.2 GzLM Structure

The GzLM consists of three elements.

1. A probability distribution from the exponential family.
2. A linear predictor $\eta = X\beta$.
3. A link function g such that $E(Y) = \mu = g^{-1}(\eta)$.

11.2.3 Link Function

Definition 1 : The link function provides the relationship between the linear predictor and the mean of the distribution function. There are many commonly used link functions, and their choice can be somewhat arbitrary. It can be convenient to match the domain of the link function to the range of the distribution function's mean.

Definition 2 : A link function is the function that links the linear model specified in the design matrix, where columns represent the beta parameters and rows the real parameters.

11.2.4 Canonical parameter

θ , called the dispersion parameter,

11.2.5 Dispersion parameter

τ , called the dispersion parameter, typically is known and is usually related to the variance of the distribution.

11.2.6 Iteratively weighted least square

IWLS is used to find the maximum likelihood estimates of a generalized linear model.

Definition: An iterative algorithm for fitting a linear model in the case where the data may contain outliers that would distort the parameter estimates if other estimation procedures were used. The procedure uses weighted least squares, the influence of an outlier being reduced by giving that observation a small weight. The weights chosen in one iteration are related to the magnitudes of the residuals in the previous iteration with a large residual earning a small weight.

11.2.7 Residual Components

In GzLMS the deviance is the sum of the deviance components

$$D = \sum d_i \quad (11.2)$$

In GzLMS the deviance is the sum of the deviance components

11.3 Generalized linear mixed models

[pawitan section 17.8]

The Generalized linear mixed model (GLMM) extend classical mixed models to non-normal outcome data.

In statistics, a generalized linear mixed model (GLMM) is a particular type of mixed model. It is an extension to the generalized linear model in which the linear predictor contains random effects in addition to the usual fixed effects. These random effects are usually assumed to have a normal distribution.

Fitting such models by maximum likelihood involves integrating over these random effects.

11.4 Assessment of Agreements in Linear and Generalized Linear Mixed Models

- Study of measuring agreement is intend to evaluate whether the readings from one rater/ measurement agree with those from other raters/measurements. In this dissertation, we are going to present a general method to assess agreement for a large variety of data with repeated measurements using linear and generalized linear mixed models.
- In the first place, a set of agreement statistics, including mean square deviation, concordance correlation coefficient, precision and accuracy coefficients, is presented for evaluating the intra-, inter-, and total-rater agreement in the multiple-rater and multiple-replications cases.
- Secondly, likelihood-based approaches are developed to estimate all the agreement statistics. Asymptotic properties of these estimates are also discussed for different data structures.
- Furthermore, our method has the merit of handling missing values and covariates naturally, and a new set of restricted agreement statistics is proposed in order to capture the true random variations and between-instrument effects adjusted for the covariate effects.
- Simulations for both linear and generalized linear mixed models are conducted to show the accuracy and effectiveness of our approaches. In the end, two industry datasets are evaluated using our approach.
- One is the cardiac function measurements used to determine the agreement between impedance cardiography and radionuclide ventriculography estimates, and the other one is an antihypertensive patch dataset given by FDA for assessing individual bioequivalence.

? generalize this approach to account for situations where the distributions are not identical, which is commonly the case. The TDI is not consistent and may not preserve its asymptotic nominal level, and that the coverage probability approach of Lin et al. (2002) is overly conservative for moderate sample sizes. This methodology proposed by ? is a regression based approach that models the mean and the variance of differences as functions of observed values of the average of the paired measurements.

Maximum likelihood estimation is used to estimate the parameters. The REML estimation is not considered since it does not lead to a joint distribution of the estimates of fixed effects and random effects parameters, upon which the assessment of agreement is based.

11.5 Random Effects and MCS

The methodology comprises two calculations. The second calculation is for the standard deviation of means. Before the modified Bland and Altman method can be applied for repeated measurement data, a check of the assumption that the variance of the repeated measurements for each subject by each method is independent of the mean of the repeated measures. This can be done by plotting the within-subject standard deviation against the mean of each subject by each method. Mean Square deviation measures the total deviation of a

11.5.1 Random coefficient growth curve model

(Chincilli 1996) Random coefficient growth curve model, a special type of mixed model have been proposed a single measure of agreement for repeated measurements.

$$\mathbf{d} = \mathbf{X}\mathbf{b} + \mathbf{Z}\mathbf{u} + \mathbf{e} \quad (11.3)$$

The distributional assumptions also require \mathbf{d} to \mathbf{N}

11.6 Random effects Model

Myles (2007) proposes the use of Random effects models to address the issue of repeated measurement.

Myles proposes a formulation of the BlandAltman plot, using the within-subject variance estimated by the random effects model, with the time of the measurement taken as a random effect. He states that *random effects models account for the dependent nature of the data, and additional explanatory variables, to provide reliable estimates of agreement in this setting.*

Agreement between methods is reflected by the between-subject variation. The Random Effects Model takes this into account before calculating the within-subject standard deviation.

11.6.1 Myers Random Effects Model

The presentation of the 95% limits of agreement is for visual judgement of how well two methods of measurement agree. The smaller the range between the two, the better the agreement is. The question of small is small is a question of clinical judgement.

Repeated measurements for each subjects are often used in clinical research.

11.6.2 Random Effects Modelling

Random effects models are used to examine the within-subject variation after adjusting for known and unknown variables, in which each subject has a different intercept and slope over a time period period.

Myles (2007) remarks that the random effects model is an extension of the analysis of variance method, accounting for more covariates.

A random effect (in Myles's case, time of measurement) is chosen to reflect the different intercept and slope for each subject with respect to their change of measurements over the time period.

In Myles's methodology, the standard deviation of difference between the means of the repeated measurements can be calculated based on the within-subject standard deviation estimates.

A random effects model (also variance components model) is a type of hierarchical linear model. Hierarchical linear modelling (HLM) is a more advanced form of simple linear regression and multiple linear regression. HLM is appropriate for use with nested data.

Faraway comments that the random effects approach is *more ambitious than the LME model in that it attempts to say something about the wider population beyond the particular sample.*

11.7 Other Approaches : Marginal Modelling

(Diggle 2002) proposes the use of marginal models as an alternative to mixed models. Marginal models are appropriate when inferences about the mean response are of specific interest.

11.8 Other Approaches

? generalize this approach to account for situations where the distributions are not identical, which is commonly the case. The TDI is not consistent and may not preserve its asymptotic nominal level, and that the coverage probability approach of Lin et al. (2002) is overly conservative for moderate sample sizes. This methodology proposed by ? is a regression based approach that models the mean and the variance of differences as functions of observed values of the average of the paired measurements.

1. Agreement and Method Comparison Studies

- (a) What is Agreement?
- (b) Repeatability
- (c)
- (d)
- (e)

2. Bland Altman Single Observations

- (a)
- (b)

3. Alternative Methods

- (a) Deming Regression
- (b) Mountain Plot
- (c) Bartko's Ellipse

- (d) Formal Tests and Procedures
- 4. Replicate Observations
- 5. LME models
- 6. Estimation and Algorithms
 - (a) ML and REML estimation
 - (b) MINQUE
 - (c)
- 7. Residual Diagnostics
 - (a) Marginal and Conditional Diagnostics
 - (b) Scaled Residuals
- 8. Influence Diagnostics
 - (a) Underlying Concepts
 - (b) Managing the Covariance Parameters
 - (c) Predicted Values, PRESS Residual and the PRESS Statistic
 - (d) Leverage
 - (e) Internally and Externally Studentized Residuals
 - (f) DFFITs and MDFFITs
 - (g) Covariance Ratio and Trace
 - (h) Likelihood Distance
 - (i) Non-iterative Update Procedures

11.9 MCS Data Sets

1. Blood Data
 2. Cardiac Data
 3. Nadler Hurley
- Introduction to Method Comparison Studies
 - Accuracy and Precision
 - Repeatability (Bland Altman 1999)
 - Barnharts Paper
 -
 - Bland and Altman Plot
 - Bland and Altman 1983 and 86
 - Limits of Agreement
 -
 -

11.10 Introduction

Outliers and detection of influent observations is an important step in the analysis of a data set. There are several ways of evaluating the influence of perturbations in the data set and in the model given the parameter estimates.

11.10.1 Overview of R implementations

Further to previous material, an appraisal of the current state of development (or lack thereof) for current implemenations for LME models, particularly for `nlme` and `lme4` fitted models.

Crucially, a review of internet resources indicates that almost all of the progress in this regard has been done for **lme4** fitted models, specifically the *Influence.ME* R package. (Nieuwenhuis et al 2014) Conversely there is very little for **nlme** models. One would immediately look at the current development workflow for both packages.

As an aside, Douglas Bates was arguably the most prominent R developer working in the LME area. However Bates has now prioritised the development of LME models in another computing environment, i.e Julia.

With regards to **nlme**, the package is now maintained by the R core development team. The most recent major text is by Galecki & Burzykowski, who have published *Linear Mixed Effects Models using R*. Also, the accompanying R package, nlmeU package is under current development, with a version being released 0.70 – 3.

The **lme4** package is used to fit linear and generalized linear mixed-effects models in the R environment. The **lme4** package is also under active development, under the leadership of Ben Bolker (McMaster Uni., Canada).

Important Consideration for MCS

The key issue is that **nlme** allows for the particular specification of Roy's Model, specifically direct specification of the VC matrices for within subject and between subject residuals. The **lme4** package does not allow for Roy's Model, for reasons that will be identified shortly. To advance the ideas that emanate from Roy's paper, one is required to use the **nlme** context. However, to take advantage of the infrastructure already provided for **lme4** models, one may change the research question away from that of Roy's paper. To this end, an exploration of what *influence.ME* can accomplish is merited.

11.11 Computation and Notation

with \mathbf{V} unknown, a standard practice for estimating $\mathbf{X}\boldsymbol{\beta}$ is to estimate the variance components σ_j^2 , compute an estimate for \mathbf{V} and then compute the projector matrix \mathbf{A} , $\mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{A}\mathbf{Y}$.

Zewotir remarks that \mathbf{D} is a block diagonal with the i -th block being $u\mathbf{I}$

11.12 Lai Shiao

? use mixed models to determine the factors that affect the difference of two methods of measurement using the conventional formulation of linear mixed effects models.

If the parameter \mathbf{b} , and the variance components are not significantly different from zero, the conclusion that there is no inter-method bias can be drawn. If the fixed effects component contains only the intercept, and a simple correlation coefficient is used, then the estimate of the intercept in the model is the inter-method bias. Conversely the estimates for the fixed effects factors can advise the respective influences each factor has on the differences. The Proc Mixed package allows users to specify different correlation structures of the variance components \mathbf{G} and \mathbf{R} .

Oxygen saturation is one of the most frequently measured variables in clinical nursing studies. ‘Fractional saturation’ (HbO_2) is considered to be the gold standard method of measurement, with ‘functional saturation’ (SO_2) being an alternative method. The method of examining the causes of differences between these two methods is applied to a clinical study conducted by ?. This experiment was conducted by 8 lab practitioners on blood samples, with varying levels of haemoglobin, from two donors. The samples have been in storage for varying periods (described by the variable ‘Bloodage’) and are categorized according to haemoglobin percentages (i.e 0%,20%,40%,60%,80%,100%). There are 625 observations in all.

? fits two models on this data, with the lab technicians and the replicate measurements as the random effects in both models. The first model uses haemoglobin level as a fixed effects component. For the second model, blood age is added as a second fixed factor.

Single fixed effect

The first model fitted by ? takes the blood level as the sole fixed effect to be analyzed. The following coefficient estimates are estimated by ‘Proc Mixed’;

$$\text{fixed effects : } 2.5056 - 0.0263\text{Fhbperct}_{ijtl} \quad (11.4)$$

$$(\text{p-values : } = 0.0054, < 0.0001, < 0.0001)$$

$$\text{random effects : } u(\sigma^2 = 3.1826) + e_{ijtl}(\sigma_e^2 = 0.1525, \rho = 0.6978)$$

$$(\text{p-values : } = 0.8113, < 0.0001, < 0.0001)$$

With the intercept estimate being both non-zero and statistically significant ($p = 0.0054$), this models supports the presence inter-method bias is 2.5% in favour of SO_2 . Also, the negative value of the haemoglobin level coefficient indicate that differences will decrease by 0.0263% for every percentage increase in the haemoglobin .

In the random effects estimates, the variance due to the practitioners is 3.1826, indicating that there is a significant variation due to technicians ($p = 0.0311$) affecting the differences. The variance for the estimates is given as 0.1525, ($p < 0.0001$).

Two fixed effects

Blood age is added as a second fixed factor to the model, whereupon new estimates are calculated;

$$\text{fixed effects : } -0.2866 + 0.1072\text{Bloodage}_{ijtl} - 0.0264\text{Fhbperct}_{ijtl}$$

$$(\text{p-values : } = 0.8113, < 0.0001, < 0.0001)$$

$$\text{random effects : } u(\sigma^2 = 10.2346) + e_{ijtl}(\sigma_e^2 = 0.0920, \rho = 0.5577)$$

$$(\text{p-values : } = 0.0446, < 0.0001, < 0.0001) \quad (11.5)$$

With this extra fixed effect added to the model, the intercept term is no longer statistically significant. Therefore, with the presence of the second fixed factor, the model is no longer supporting the

presence of inter-method bias. Furthermore, the second coefficient indicates that the blood age of the observation has a significant bearing on the size of the difference between both methods ($p < 0.0001$). Longer storage times for blood will lead to higher levels of particular blood factors such as MetHb and HbCO (due to the breakdown and oxidisation of the haemoglobin). Increased levels of MetHb and HbCO are concluded to be the cause of the differences. The coefficient for the haemoglobin level doesn't differ greatly from the single fixed factor model, and has a much smaller effect on the differences. The random effects estimates also indicate significant variation for the various technicians; 10.2346 with $p = 0.0446$.

? demonstrates how that linear mixed effects models can be used to provide greater insight into the cause of the differences. Naturally the addition of further factors to the model provides for more insight into the behavior of the data.

11.13 Liao Shaio

Lai et Shiao is interesting in that it extends the usual method comparison study question. It correctly identifies LME models as a methodology that can be used to make such questions tractable. The Data Set used in their examples is unavailable for independent use. Therefore, for the sake of consistency, a data set will be simulated based on the Blood Data that will allow for extra variables.

A Study of the Bland-Altman Plot and its Associated Methodology

Joseph G. Voelkel Bruce E. Siskowski

11.14 Limits of agreement for Carstensen's data

? describes the calculation of the limits of agreement (with the inter-method bias implicit) for both data sets, based on his formulation;

$$\hat{\alpha}_1 - \hat{\alpha}_2 \pm 2\sqrt{2\hat{\tau}^2 + \hat{\sigma}_1^2 + \hat{\sigma}_2^2}.$$

For the 'Fat' data set, the inter-method bias is shown to be 0.045. The limits of agreement are $(-0.23, 0.32)$

Carstensen demonstrates the use of the interaction term when computing the limits of agreement for the 'Oximetry' data set. When the interaction term is omitted, the limits of agreement are $(-9.97, 14.81)$. Carstensen advises the inclusion of the interaction term for linked replicates, and hence the limits of agreement are recomputed as $(-12.18, 17.12)$.

11.15 Hamlett and Lam

The methodology proposed by ? is largely based on Hamlett et al. (2004), which in turn follows on from Lam et al. (1999).

The desired outcome of this research is to

- Formulate a methodology that represents Best practice in Method Comparison Studies. Indeed the methodology is envisaged to advance what is considered best practice, inter alia, by making diagnostics procedures a standard part of MCS.
- Provide for ease of use such that non-statisticians can master and implement the method, with a level of training that one would expect as part of a Professional CPD programme.

Apropos of the matter of ease-of-use, certain assumptions must be made.

The user has a reasonable amount of computer literacy. The user would have a reasonable understanding of statistics, consistent with an undergraduate statistics module. That is to say, that the user is acquainted with the idea of p -values.

Easy to follow set of instructions to properly implement the method.

Linear Mixed Effects Models can be implemented by using one of the following R packages. lme4
nlme

The first package to be introduced was nlme, developed by Jose Pinheiro and Douglas Bates (Authors of the the companion textbook, NAME)

As this package has been under ongoing development for quite a long time, it is now allows for a lot of complex LME implementations. Furthermore, nlme is one of the base R packages. That is to say, when one downloads and installs R, nlme is automatically installed also, and can be called immediately.

Having said that, the authors have pointed to several limitations of the overall methodology through R. The original developers have both left the project, but other statisticians have taken over the development, and indeed a new version of nlme was released.

LME4 is a more recent package. at a glance, the syntax is easier, but the development is less advanced. There are several functionalities that can not be implemented with lme4 yet. As an example - CHAP5 in PB - has no equivalent in LME4. Indeed no textbook exists to co-incide with LME4.

The main author, Douglas Bates, has turned his attention to development of LME models in the Julia programming language.

The nlmeU package is described by its authors as an extension of the nlme package, and indeed provides for additional functionality. The package is also useful as it serves as a companion piece to the book by Galecki and Burzwhatski.

The nlme package also allows for the specification of GLS models.

Objects and Classes

The main nlme object is an `nlme` model.

The main lme4 object is called an `lmer` model

The lattice package is used for graphical methods.

Model Diagnostics with `nlme`

11.15.1 Inappropriate Techniques for MCS

11.15.2 Links and Papers

Westgard Statistics - <http://www.westgard.com/lesson23.htm>

Measurement Systems Analysis

The topic of measurement sensitivity analysis (MSA, also known as Gauge R&R) is prevalent in industrial statistics (i.e Six Sigma).

There is extensive literature that covers the area. For the sake of brevity, we will use Cano et al.

For sake of clarity, Cano's definitions of repeatability and reproducibility are listed, with added emphasis.

Reproducibility is rarely, if ever, discussed in the domain of Method Comparison Studies. This may be due to the fact that prevalent methodologies can be used for the problem. However the methodologies proposed by this research can easily be extended.

Chapter 12

Introduction

Chapter 13

Appendix

Bayesian BA - Philip J Schluter

Bayesian Bland Altman Approaches A multivariate hierarchical Bayesian approach to measuring agreement in repeated measurement method comparison studies

<http://www.biomedcentral.com/1471-2288/9/6>

Background

Assessing agreement in method comparison studies depends on two fundamentally important components; validity (the between method agreement) and reproducibility (the within method agreement).

The Bland-Altman limits of agreement technique is one of the favoured approaches in medical literature for assessing between method validity. However, few researchers have adopted this approach for the assessment of both validity and reproducibility.

This may be partly due to a lack of a flexible, easily implemented and readily available statistical machinery to analyse repeated measurement method comparison data.

Methods

Adopting the Bland-Altman framework, but using Bayesian methods, we present this statistical machinery. Two multivariate hierarchical Bayesian models are advocated, one which assumes that the underlying values for subjects remain static (exchangeable replicates) and one which assumes that the

underlying values can change between repeated measurements (non-exchangeable replicates).

Results

We illustrate the salient advantages of these models using two separate datasets that have been previously analysed and presented; (i) assuming static underlying values analysed using both multivariate hierarchical Bayesian models, (ii) assuming each subject's underlying value is continually changing quantity and analysed using the non-exchangeable replicate multivariate hierarchical Bayesian model.

Conclusion These easily implemented models allow for full parameter uncertainty, simultaneous method comparison, handle unbalanced or missing data, and provide estimates and credible regions for all the parameters of interest. Computer code for the analyses is also presented, provided in the freely available and currently cost free software package WinBUGS.

Bayesian Approach

A multivariate hierarchical Bayesian approach to measuring agreement in repeated measurement method comparison studies PJ Schluter - BMC medical research methodology, 2009 - biomedcentral.com

- Assessing agreement in method comparison studies depends on two fundamentally important components; validity (the between method agreement) and reproducibility (the within method agreement).
- The Bland-Altman limits of agreement technique is one of the f

13.1 Escaramis

13.1.1 Background

In an agreement assay, it is of interest to evaluate the degree of agreement between the different methods (devices, instruments or observers) used to measure the same characteristic. We propose in this study a technical simplification for inference about the total deviation index (TDI) estimate to assess agreement between two devices of normally-distributed measurements and describe its utility to evaluate inter- and intra-rater agreement if more than one reading per subject is available for each device.

13.1.2 Methods

We propose to estimate the TDI by constructing a probability interval of the difference in paired measurements between devices, and thereafter, we derive a tolerance interval (TI) procedure as a natural way to make inferences about probability limit estimates. We also describe how the proposed method can be used to compute bounds of the coverage probability.

13.1.3 Results

The approach is illustrated in a real case example where the agreement between two instruments, a handle mercury sphygmomanometer device and an OMRON 711 automatic device, is assessed in a sample of 384 subjects where measures of systolic blood pressure were taken twice by each device. A simulation study procedure is implemented to evaluate and compare the accuracy of the approach to two already established methods, showing that the TI approximation produces accurate empirical confidence levels which are reasonably close to the nominal confidence level.

13.1.4 Conclusions

The method proposed is straightforward since the TDI estimate is derived directly from a probability interval of a normally-distributed variable in its original scale, without further transformations. Thereafter, a natural way of making inferences about this estimate is to derive the appropriate TI. Constructions of TI based on normal populations are implemented in most standard statistical packages, thus making it simpler for any practitioner to implement our proposal to assess agreement.

Lin defined the TDI as the boundary, κ_P which captures a large proportion p of paired based differences from two devices or observers within the boundary.

The value of κ_P that yields $P(|D| < \kappa_P) = p$ where D is the paired-difference variate.

$$\kappa_P = F^{-1}(p) = \sigma_D \sqrt{\chi^2(p, 1, \mu_D^2/\sigma_d^2)}$$

$$\kappa_P = Z_{\frac{1+p}{2}} \|\varepsilon\|$$

$$\hat{\kappa}_p = \hat{\mu}_D = Z_{p_i} \sigma_d$$

Coverage Probability is another user friendly measure of agreement which is related to the computation of the TDI.

13.2 Schabenberger

schab examines the use and implementation of influence measures in LME models.

Influence is understood to be the ability of a single or multiple data points, through their presences or absence in the data, to alter important aspects of the analysis, yield qualitatively different inferences, or violate assumptions of the statistical model (*schabenberger*).

Outliers are the most noteworthy data points in an analysis, and an objective of influence analysis is how influential they are, and the manner in which they are influential.

schab describes a simple procedure for quantifying influence. Firstly a model should be fitted to the data, and estimates of the parameters should be obtained. The second step is that either single or multiple data points, specifically outliers, should be omitted from the analysis, with the original parameter estimates being updated.

This is known as ‘*leave one out*’ *leave k out*’ analysis. The final step of the procedure is comparing the sets of estimates computed from the entire and reduced data sets to determine whether the absence of observations changed the analysis.

schabenberger notes that it is not always possible to derive influence statistics necessary for comparing full- and reduced-data parameter estimates.

In recent years, mixed models have become invaluable tools in the analysis of experimental and observational data. In these models, more than one term can be subject to random variation. Mixed model technology enables you to analyze complex experimental data with hierarchical random processes, temporal, longitudinal, and spatial data, to name just a few important applications.

schab remarks that the concept of critiquing the model-data agreement applies in mixed models in the same way as in linear fixed-effects models. In fact, because of the more complex model structure,

you can argue that model and data diagnostics are even more important. For example, you are not only concerned with capturing the important variables in the model. You are also concerned with “distributing them correctly between the fixed and random components of the model. The mixed model structure presents unique and interesting challenges that prompt us to reexamine the traditional ideas of influence and residual analysis.

13.3 Hawkins : Diagnostics for conformity of paired quantitative measurements

- Matched pairs data arise in many contexts in case-control clinical trials, for example, and from cross-over designs. They also arise in experiments to verify the equivalence of quantitative assays. This latter use (which is the main focus of this paper) raises difficulties not always seen in other matched pairs applications.
- Since the designs deliberately vary the analyte levels over a wide range, issues of variance dependent on mean, calibrations of differing slopes, and curvature all need to be added to the usual model assumptions such as normality.
- Violations in any of these assumptions invalidate the conventional matched pairs analysis.
- A graphical method, due to Bland and Altman, of looking at the relationship between the average and the difference of the members of the pairs is shown to correspond to a formal testable regression model.
- Using standard regression diagnostics, one may detect and diagnose departures from the model assumptions and remedy them for example using variable transformations. Examples of different common scenarios and possible approaches to handling them are shown.

A multi-Rate nonparametric test of agreement and corresponding agreement plot

- Published in: Computational Statistics and Data Analysis 54(2010)109-119 - Author: Alan D. Hutson, University of Buffalo

This approach takes advantage of readily available tests of uniformity found in most statistical software packages. Such tests include the KS d statistic, the Anderson Darling Statistic and the Cramer-Von Mises statistical test for univariate data.

An important aspect of this approach is the "Agreement Region".

Roy Test

Roy's Tests (Roy 2009) Roy 2009 devised an LME based Testing approach to the MCS problem, based on earlier work by Hamlett et al. Roy 2009 presents a series of three formal hypothesis tests for assessing agreement between two methods of measurement. Roy also alludes to some of the current shortcomings of the approach.

Comparing different model specifications with LRT tests

- Roy 2007 - Roy 2009 - Hamlett et al. - Roy Leiva 2011

Conventionally LME models can be tested using Likelihood Ratio Tests, wherein a reference model is compared to a nested model.

```
> Ref.Fit = lme(y ~ meth-1, data = dat, #Symm , Symm#
+ random = list(item=pdSymm(~ meth-1)),
+ weights=varIdent(form=~1|meth),
+ correlation = corSymm(form=~1 | item/repl),
+ method="ML")
```

Roy(2009) presents two nested models that specify the condition of equality as required, with a third nested model for an additional test. These three formulations share the same structure, and can be specified by making slight alterations of the code for the Reference Model.

Nested Model (Between-Item Variability)

```
> NMB.fit = lme(y ~ meth-1, data = dat, #CS , Symm#
+ random = list(item=pdCompSymm(~ meth-1)),
+ correlation = corSymm(form=~1 | item/repl),
+ method="ML")
```

Nested Model (Within item Variability)

```
> NMW.fit = lme(y ~ meth-1, data = dat,    #Symm , CS#
+   random = list(item=pdSymm(~ meth-1)),
+   weights=varIdent(form=~1|meth),
+   correlation = corCompSymm(form=~1 | item/repl),
+   method="ML")
```

Nested Model (Overall Variability) Additionally there is a third nested model, that can be used to test overall variability, substantively a a joint test for between-item and within-item variability. The motivation for including such a test in the suite is not clear, although it does circumvent the need for multiple comparison procedures in certain circumstances, hence providing a simplified procedure for non-statisticians.

```
> NMO.fit = lme(y ~ meth-1, data = dat,    #CS , CS#
+   random = list(item=pdCompSymm(~ meth-1)),
+   correlation = corCompSymm(form=~1 | item/repl),
+   method="ML")
```

ANOVAs for Original Fits The likelihood Ratio test is very simple to implement in R. All that is required it to specify the reference model and the relevant nested mode as arguments to the command `anova()`. The figure below displays the three tests described by Roy (2009).

```
> testB    = anova(Ref.Fit,NMB.fit)                # Between-Subject Variabilities
> testW    = anova(Ref.Fit,NMW.fit)                # Within-Subject Variabilities
> testO    = anova(Ref.Fit,NMO.fit)                # Overall Variabilities
```

13.4 Profile Function with "lmer"

The `profile()` function for lmer models is now available in the latest version of lme4, to be installed by typing:

```
install.packages("lme4",repos="http://r-forge.r-project.org")  
also
```

The `mle` function from the `stats4` package is a wrapper of `optim`, which makes it quite easy to produce profile likelihood computations.

See `help("profile,mle-method", package = "stats4")` for more information.

<http://people.upei.ca/hstryhn/stryhn208.pdf>

The profile likelihood (or likelihood or likelihood ratio) method is applicable to all likelihood based statistical analysis and is generally less sensitive to the difficulties encountered by Wald-Tyoe CIs.

13.5 Turkan's LMEs

The linear mixed model is considerably sensitive to outliers and influential observations. It is known that outliers and influential observations affect substantially the results of analysis. So it is very important to be aware of these observations.

Some diagnostics which are analogue of diagnostics in multiple linear regression were developed to detect outliers and influential observations in the linear mixed model. *In this paper, the new diagnostic measure which is analogue of the Pena's influence statistic is developed for the linear mixed model.*

$$\hat{u} = DZ^T H^{-1}(y - X\hat{\beta})$$

$$\hat{y} = (I_n - H^{-1})y + H^{-1}X\hat{\beta}$$

The proposed diagnostic Measure.

13.5.1 Ordinary Least Product Regression

Ludbrook (1997) states that the grouping structure can be straightforward, but there are more complex data sets that have a hierarchical(nested) model.

Observations between groups are independent, but observations within each groups are dependent because they belong to the same subpopulation. Therefore there are two sources of variation: between-group and within-group variance. Mean correction is a method of reducing bias.

13.5.2 A regression based approach based on Bland Altman Analysis

Lu et al used such a technique in their comparison of DXA scanners. They also used the Blackwood Bradley test. However it was shown that, for particular comparisons, agreement between methods was indicated according to one test, but lack of agreement was indicated by the other.

13.6 Measurement Error Models

Dunn (2002) proposes a measurement error model for use in method comparison studies. Consider n pairs of measurements X_i and Y_i for $i = 1, 2, \dots, n$.

$$X_i = \tau_i + \delta_i \tag{13.1}$$

$$Y_i = \alpha + \beta\tau_i + \epsilon_i$$

In the above formulation is in the form of a linear structural relationship, with τ_i and $\beta\tau_i$ as the true values, and δ_i and ϵ_i as the corresponding measurement errors. In the case where the units of measurement are the same, then $\beta = 1$.

$$E(X_i) = \tau_i \quad (13.2)$$

$$E(Y_i) = \alpha + \beta\tau_i$$

$$E(\delta_i) = E(\epsilon_i) = 0$$

The value α is the inter-method bias between the two methods.

$$z_0 = d = 0 \quad (13.3)$$

$$z_{n+1} = z_n^2 + c \quad (13.4)$$

13.7 Work List

1. ML v REML
2. Nested Models and LRTs
3. Generalized Least Squares
4. Diagnostics
5. Simplifying GLS
6. Paper progression

13.8 Diagnostics

13.8.1 Identifying outliers with a LME model object

The process is slightly different than with standard LME model objects, since the *influence* function does not work on lme model objects. Given *mod.lme*, we can use the plot function to identify outliers.

13.8.2 Diagnostics for Random Effects

Empirical best linear unbiased predictors EBLUPS provide the a useful way of diagnosing random effects.

EBLUPs are also known as “shrinkage estimators” because they tend to be smaller than the estimated effects would be if they were computed by treating a random factor as if it was fixed (West et al)

13.9 Two-tailed testing

A test for equality of variances, based on the likelihood Ratio test, is very simple to implement using existing methodologies. All that is required is to specify the reference model and the relevant nested mode as arguments to the command `anova()`. The output can be interpreted in the usual way.

13.10 One Tailed Testing

The approach proposed by Roy deals with the question of agreement, and indeed interchangeability, as developed by Bland and Altman's corpus of work. In the view of Dunn, a question relevant to many practitioners is which of the two methods is more precise.

The relationship between precision and the within-item and between-item variability must be established. Roy establishes the equivalence of repeatability and within-item variability, and hence precision. The method with the smaller within-item variability can be deemed to be the more precise.

13.11 Enabling One Tailed Testing

A useful approach is to compute the confidence intervals for the ratio of within-item standard deviations (equivalent to the ratio of repeatability coefficients), which can be interpreted in the usual manner (or alternatively, the ratio of the variances). In fact, the ratio of within-item standard deviations, with the attendant confidence interval, can be determined using a single R command: `intervals()`.

Pinheiro and Bates (pg 93-95) give a description of how confidence intervals for the variance components are computed. Furthermore a complete set of confidence intervals can be computed to complement the variance component estimates. However, to facilitate one tailed testing, what is required is the computation of the variance ratios of within-item and between-item standard deviations.

A naive approach would be to compute the variance ratios by relevant F distribution quantiles. However, the question arises as to the appropriate degrees of freedom. However, Douglas Bates has stated that an alternative approach is required (i.e. Profile Likelihoods)

"The omission of standard errors on variance components is intentional. The distribution

of an estimator of a variance component is highly skewed and obtaining an estimate of the standard deviation of a skewed distribution is not very useful. A much better approach is based on profiling the objective function.” (Douglas Bates May 2012)

13.12 Profile Likelihood

Normal-based confidence intervals for a parameter of interest are inaccurate when the sampling distribution of the estimate is skewed. The technique known as profile likelihood can produce confidence intervals with better coverage. It may be used when the model includes only the variable of interest or several other variables in addition. Profile-likelihood confidence intervals are particularly useful in nonlinear models.

Profile likelihood confidence intervals are based on the log-likelihood function.

13.13 Implementation of PL Confidence Intervals

The suitable calculation of confidence limits for this variance ratio are to be computed using the profile likelihood approach. The R package `profilelikelihood` will be assessed for feasibility, particularly the command `profilelikelihood.lme()`

13.14 Extension of techniques to LME Models

Model diagnostic techniques, well established for classical models, have since been adapted for use with linear mixed effects models. Diagnostic techniques for LME models are inevitably more difficult to implement, due to the increased complexity.

Beckman, Nachtsheim and Cook (1987) ? applied the local influence method of Cook (1986) to the analysis of the linear mixed model.

While the concept of influence analysis is straightforward, implementation in mixed models is more complex. Update formulae for fixed effects models are available only when the covariance parameters are assumed to be known.

If the global measure suggests that the points in U are influential, the nature of that influence should be determined. In particular, the points in U can affect the following

- the estimates of fixed effects,
- the estimates of the precision of the fixed effects,
- the estimates of the covariance parameters,
- the estimates of the precision of the covariance parameters,
- fitted and predicted values.

13.15 Covariance Parameters

The unknown variance elements are referred to as the covariance parameters and collected in the vector θ .

13.15.1 Methods and Measures

The key to making deletion diagnostics useable is the development of efficient computational formulas, allowing one to obtain the case deletion diagnostics by making use of basic building blocks, computed only once for the full model.

Zewotir and Galpin (2005) lists several established methods of analyzing influence in LME models. These methods include

- Cook’s distance for LME models,
- likelihood distance,
- the variance (information) ration,
- the Cook-Weisberg statistic,
- the Andrews-Prebigon statistic.

13.16 Computation and Notation

with \mathbf{V} unknown, a standard practice for estimating $\mathbf{X}\boldsymbol{\beta}$ is the estimate the variance components σ_j^2 , compute an estimate for \mathbf{V} and then compute the projector matrix A , $\mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{A}\mathbf{Y}$.

Zewotir and Galpin (2005) remarks that \mathbf{D} is a block diagonal with the i -th block being $u\mathbf{I}$

13.17 Measures 2

13.17.1 Cook’s Distance

- For variance components γ

Diagnostic tool for variance components

$$C_{\theta i} = ((\hat{\theta})_{[i]} - \hat{\theta})^T \text{cov}(\hat{\theta})^{-1} ((\hat{\theta})_{[i]} - \hat{\theta})$$

13.18 Carstensen’s Mixed Models

Carstensen (2004) proposes linear mixed effects models for deriving conversion calculations similar to Deming’s regression, and for estimating variance components for measurements by different methods.

The following model (in the authors own notation) is formulated as follows, where y_{mir} is the r th replicate measurement on subject i with method m .

$$y_{mir} = \alpha_m + \beta_m \mu_i + c_{mi} + e_{mir} \quad (e_{mi} \sim N(0, \sigma_m^2), c_{mi} \sim N(0, \tau_m^2)) \quad (13.5)$$

The intercept term α and the $\beta_m \mu_i$ term follow from Dunn (2002), expressing constant and proportional bias respectively , in the presence of a real value μ_i . c_{mi} is a interaction term to account for replicate, and e_{mir} is the residual associated with each observation. Since variances are specific to each method, this model can be fitted separately for each method.

The above formulation doesn't require the data set to be balanced. However, it does require a sufficient large number of replicates and measurements to overcome the problem of identifiability. The import of which is that more than two methods of measurement may be required to carry out the analysis. There is also the assumptions that observations of measurements by particular methods are exchangeable within subjects. (Exchangeability means that future samples from a population behaves like earlier samples).

Carstensen (2004) uses the above formula to predict observations for a specific individual i by method m ;

$$BLUP_{mir} = \hat{\alpha}_m + \hat{\beta}_m \mu_i + c_{mi} \quad (13.6)$$

. Under the assumption that the μ s are the true item values, this would be sufficient to estimate parameters. When that assumption doesn't hold, regression techniques (known as updating techniques) can be used additionally to determine the estimates. The assumption of exchangeability can be unrealistic in certain situations. Carstensen (2004) provides an amended formulation which includes an extra interaction term ($d_{mr} \sim N(0, \omega_m^2)$)to account for this.

Carstensen et al. (2008) sets out a methodology of computing the limits of agreement based upon variance component estimates derived using linear mixed effects models. Measures of repeatability, a characteristic of individual methods of measurements, are also derived using this method.

13.18.1 Using LME models to create Prediction Intervals

Carstensen (2004) also advocates the use of linear mixed models in the study of method comparisons. The model is constructed to describe the relationship between a value of measurement and its real value. The non-replicate case is considered first, as it is the context of the Bland-Altman plots. This model assumes that inter-method bias is the only difference between the two methods. A measurement y_{mi} by method m on individual i is formulated as follows;

$$y_{mi} = \alpha_m + \mu_i + e_{mi} \quad (e_{mi} \sim N(0, \sigma_m^2)) \quad (13.7)$$

The differences are expressed as $d_i = y_{1i} - y_{2i}$. For the replicate case, an interaction term c is added to the model, with an associated variance component. All the random effects are assumed independent, and that all replicate measurements are assumed to be exchangeable within each method.

$$y_{mir} = \alpha_m + \mu_i + c_{mi} + e_{mir} \quad (e_{mi} \sim N(0, \sigma_m^2), c_{mi} \sim N(0, \tau_m^2)) \quad (13.8)$$

Chapter 14

Model Diagnostics

Contents

Abstract

This chapter is broken into two parts. The first part is a review of diagnostics methods for linear models, intended to acquaint the reader with the subject, and also to provide a basis for material covered in the second part. Particular attention is drawn to graphical methods.

The second part of the chapter looks at diagnostics techniques for LME models, firstly covering the theory, then proceeding to a discussion on implementing these using **R** code.

While a substantial body of work has been developed in this area, there is still areas worth exploring. In particular the development of graphical techniques pertinent to LME models should be looked at.

14.1 Model Validation Framework

In statistical modelling, the process of model validation is a critical step of model fitting process, but also a step that is too often overlooked. A very simple procedure is to examine commonly-used metrics, such as the R^2 value. However, using a small handful of simple measures and methods is insufficient to properly assess the quality of a fitted model. To do so properly, a full and comprehensive analysis that tests of all of the assumptions, as far as possible, must be carried out.

? describes the model validation framework as comprised of the following tasks

- overall measures of goodness-of-fit
- the informal, graphical examination of estimates of model errors to assess the quality of distributional assumptions: residual analysis
- the quantitative assessment of the inter-relationship of model components; for example, collinearity diagnostics
- the qualitative and quantitative assessment of influence of cases on the analysis, i.e. influence analysis.

Residual analysis is a widely used model validation technique. A residual is simply the difference between an observed value and the corresponding fitted value, as predicted by the model. The rationale is that, if the model is properly fitted to the model, then the residuals would approximate the random errors that one should expect. that is to say, if the residuals behave randomly, with no discernible trend, the model has fitted the data well. If some sort of non-random trend is evident in the model, then the model can be considered to be poorly fitted.

Statistical software environments, such as the R Programming language, provides a suite of tests and graphical procedure sfor appraising a fitted linear model, with several of these procedures analysing the model residuals.

For LME models the matter of residual is more complex. ? describes two types of residuals, marginal residuals and conditional residuals. A marginal residual is the difference between the observed data and the estimated marginal mean. A conditional residual is the difference between the observed data

and the predicted value of the observation. In a model without random effects, both sets of residuals coincide. We shall revert to this matter in due course.

Further to the analysis of residuals, ? recommends the examination of the following questions.

- Does the model-data agreement support the model assumptions?
- Should model components be refined, and if so, which components? For example, should regressors be added or removed, and is the covariation of the observations modeled properly?
- Are the results sensitive to model and/or data? Are individual data points or groups of cases particularly influential on the analysis?

14.1.1 Outliers and Leverage

The question of whether or not a point should be considered an outlier must also be addressed. An outlier is an observation whose true value is unusual given its value on the predictor variables. The leverage of an observation is a further consideration. Leverage describes an observation with an extreme value on a predictor variable is a point with high leverage. High leverage points can have a great amount of effect on the estimate of regression coefficients.

Influence can be thought of as the product of leverage and outlierness. An observation is said to be influential if removing the observation substantially changes the estimate of the regression coefficients. The R programming language has a variety of methods used to study each of the aspects for a linear model. While linear models and GLMS can be studied with a wide range of well-established diagnostic techniques, the choice of methodology is much more restricted for the case of LMEs.

For classical linear models, residual diagnostics are typically conducted using a plot of the observed residuals and the predicted values. A visual inspection for the presence of trends inform the analyst on the validity of distributional assumptions, and to detect outliers and influential observations.

Cook's distance

In the study of Linear model diagnostics, Cook proposed a measure that combines the information of leverage and residual of the observation, now known simply as the Cook's Distance. ? would later

adapt the Cook's distance measure for the analysis of LME models.

14.2 Zewotir Measures of Influence in LME Models

Zewotir and Galpin (2005) describes a number of approaches to model diagnostics, investigating each of the following;

- Variance components
- Fixed effects parameters
- Prediction of the response variable and of random effects
- likelihood function

Zewotir and Galpin (2005) lists several established methods of analyzing influence in LME models. These methods include

- Cook's distance for LME models,
- likelihood distance,
- the variance (information) ration,
- the Cook-Weisberg statistic,
- the Andrews-Prebigon statistic.

14.3 Matrix Notation for Case Deletion

14.3.1 Case deletion notation

For notational simplicity, $\mathbf{A}(i)$ denotes an $n \times m$ matrix \mathbf{A} with the i -th row removed, a_i denotes the i -th row of \mathbf{A} , and a_{ij} denotes the (i, j) -th element of \mathbf{A} .

14.3.2 Further Assumptions of Linear Models

As with fitted models, the assumption of normality of residuals and homogeneity of variance is applicable to LMEs also.

Homoscedascity is the technical term to describe the variance of the residuals being constant across the range of predicted values. Heteroscedascity is the converse scenario : the variance differs along the range of values.

On occasion, quantification is not possible. Assume, for example, that a data point is removed and the new estimate of the G matrix is not positive definite. This may occur if a variance component estimate now falls on the boundary of the parameter space. Thus, it may not be possible to compute certain influence statistics comparing the full-data and reduced-data parameter estimates. However, knowing that a new singularity was encountered is important qualitative information about the data points influence on the analysis.

The basic procedure for quantifying influence is simple:

1. Fit the model to the data and obtain estimates of all parameters.
2. Remove one or more data points from the analysis and compute updated estimates of model parameters.
3. Based on full- and reduced-data estimates, contrast quantities of interest to determine how the absence of the observations changes the analysis.

We use the subscript (U) to denote quantities obtained without the observations in the set U. For example, (U) denotes the fixed-effects **leave-U-out** estimates. Note that the set U can contain multiple observations.

If the global measure suggests that the points in U are influential, you should next determine the nature of that influence. In particular, the points can affect

- the estimates of fixed effects
- the estimates of the precision of the fixed effects
- the estimates of the covariance parameters
- the estimates of the precision of the covariance parameters
- fitted and predicted values

It is important to further decompose the initial finding to determine whether data points are actually troublesome. Simply because they are influential somehow, should not trigger their removal from the analysis or a change in the model. For example, if points primarily affect the precision of the covariance parameters without exerting much influence on the fixed effects, then their presence in the data may not distort hypothesis tests or confidence intervals about β .

14.3.3 Summary of Paper

Standard residual and influence diagnostics for linear models can be extended to LME models. The dependence of the fixed effects solutions on the covariance parameters has important ramifications on the perturbation analysis. Calculating the studentized residuals-And influence statistics whereas each software procedure can calculate both conditional and marginal raw residuals, only SAS Proc Mixed is currently the only program that provide studentized residuals Which are preferred for model diagnostics. The conditional Raw residuals are not well suited to detecting outliers as are the studentized conditional residuals. (schabenbege r)

LME are flexible tools for the analysis of clustered and repeated measurement data. LME extend the capabilities of standard linear models by allowing unbalanced and missing data, as long as the missing data are MAR. Structured covariance matrices for both the random effects G and the residuals R . missing at Random.

A conditional residual is the difference between the observed value and the predicted value of a dependent variable. Influence diagnostics are formal techniques that allow the identification of observations that heavily influence estimates of parameters. To alleviate the problems with the interpretation of conditional residuals that may have unequal variances, we consider scaling. Residuals obtained in this manner are called studentized residuals.

14.4 Schabenberger: Summary and Conclusions

- Standard residual and influence diagnostics for linear models can be extended to linear mixed models. The dependence of fixed-effects solutions on the covariance parameter estimates has important ramifications in perturbation analysis.
- To gauge the full impact of a set of observations on the analysis, covariance parameters need to be updated, which requires refitting of the model.
- The experimental INFLUENCE option of the MODEL statement in the MIXED procedure (SAS 9.1) enables you to perform iterative and noniterative influence analysis for individual observations and sets of observations.
- The conditional (subject-specific) and marginal (population-averaged) formulations in the linear mixed model enable you to consider conditional residuals that use the estimated BLUPs of the random effects, and marginal residuals which are deviations from the overall mean.
- Residuals using the BLUPs are useful to diagnose whether the random effects components in the model are specified correctly, marginal residuals are useful to diagnose the fixed-effects components.
- Both types of residuals are available in SAS 9.1 as an experimental option of the MODEL statement in the MIXED procedure.
- It is important to note that influence analyses are performed under the assumption that the chosen model is correct. Changing the model structure can alter the conclusions. Many other variance models have been fit to the data presented in the repeated measures example. You need to see the conclusions about which model component is affected in light of the model being fit.

Leave-One-Out Diagnostics with lmeU

Galecki et al provide a brief the matter of LME influence diagnostics in their book.

The command `lmeU` fits a model with a particular subject removed. The identifier of the subject to be removed is passed as the only argument

A plot of the per-observation diagnostics individual subject log-likelihood contributions can be rendered.

The addition of an extra factor

Interaction terms are featured in ANOVA designs.

My search just now found no mention of Cook's distance or influence measures.

The closest I found was an unanswered question on this from April 2003 (<http://finzi.psych.upenn.edu/R/Rhe>

Beyond that, there is an excellent discussion of "Examining a Fitted Model" in Sec. 4.3 (pp. 174-197) of Pinheiro and Bates (2000) *Mixed-Effects Models in S and S-Plus* (Springer).

Pinheiro and Bates decided NOT to include plots of Cook's distance among the many diagnostics they did provide. However, `plot(fit.lme)` plots 'standardized residuals' vs. predicted or 'fitted values'. Wouldn't points with large influence stand apart from the crowd in terms of 'fitted value'?

Of course, there are many things other one could do to get at related information, including reading the code for 'influence' and 'lme', and figure out from that how to write an 'influence' method for an 'lme' object.

14.5 Paired T tests

This method can be applied to test for statistically significant deviations in bias. This method can be potentially misused for method comparison studies.

It is a poor measure of agreement when the rater's measurements are perpendicular to the line of equality [Hutson et al]. In this context, an average difference of zero between the two raters, yet the scatter plot displays strong negative correlation.

Components in assessing agreement

1. The degree of linear relationship between the two sets
2. The amount of bias as represented by the difference in the means
3. The Differences in the two variances.

14.6 Methods of assessing agreement

1. Pearson's Correlation Coefficient
2. Intraclass correlation coefficient
3. Bland Altman Plot
4. Bartko's Ellipse (1994)
5. Blackwood Bradley Test
6. Lin's Reproducibility Index
7. Luiz Step function

Bland and Altman attend to the issue of repeated measures in 1996.

Repeated measurements on several subjects can be used to quantify measurement error, the variation between measurements of the same quantity on the same individual.

Bland and Altman discuss two metrics for measurement error; the within-subject standard deviation ,and the correlation coefficient.

The above plot incorporates both the conventional limits of agreement (the inner pair of dashed lines), the 't' limits of agreement (the outer pair of dashed lines) centred around the inter-method bias (indicated by the full line). This plot is intended for expository purposes only, as the sample size is small.

14.6.1 Equivalence and Interchangeability

Limits of agreement are intended to analyse equivalence. How this is assessed is the considered judgement of the practitioner. In Bland and Altman (1986) an example of good agreement is cited. For two methods of measuring ‘oxygen saturation’, the limits of agreement are calculated as $(-2.0, 2.8)$. A practitioner would ostensibly find this to be sufficiently narrow.

If the limits of agreement are not clinically important, which is to say that the differences tend not to be substantial, the two methods may be used interchangeably. Dunn (2002) takes issue with the notion of ‘equivalence’, remarking that while agreement indicated equivalence, equivalence does not reflect agreement.

Chapter 15

Model Diagnostics

15.1 Introduction

In classical linear models model diagnostics have been become a required part of any statistical analysis, and the methods are commonly available in statistical packages and standard textbooks on applied regression. However it has been noted by several papers that model diagnostics do not often accompany LME model analyses.

15.1.1 Checking model assumptions

In classical linear regression, it is important to carry out model diagnostic techniques to determine whether or not the distributional assumptions are satisfied. Model diagnostics are also used to determine the influence of unusual observations.

Schabenberger (2004) describes the examination of model-data agreement as comprising several elements; residual analysis, goodness of fit, collinearity diagnostics and influence analysis.

15.1.2 Influence Diagnostics: Basic Idea and Statistics

The general idea of quantifying the influence of one or more observations relies on computing parameter estimates based on all data points, removing the cases in question from the data, refitting the model, and computing statistics based on the change between full-data and reduced-data estimation. Influence

statistics can be coarsely grouped by the aspect of estimation that is their primary target:

- overall measures compare changes in objective functions: (restricted) likelihood distance (Cook and Weisberg 1982, Ch. 5.2)
- influence on parameter estimates: Cook's (Cook 1977, 1979), MDFFITS (Belsley, Kuh, and Welsch 1980, p. 32)
- influence on precision of estimates: CovRatio and CovTrace
- influence on fitted and predicted values: PRESS residual, PRESS statistic (Allen 1974), DFFITS (Belsley, Kuh, and Welsch 1980, p. 15)
- outlier properties: internally and externally studentized residuals, leverage

15.1.3 Introduction

The linear mixed effects model is a useful methodology for fitting a wide range of models. However, linear mixed effects models are known to be sensitive to outliers. ? advises that identification of outliers is necessary before conclusions may be drawn from the fitted model.

Standard statistical packages concentrate on calculating and testing parameter estimates without considering the diagnostics of the model.

The assessment of the effects of perturbations in data, on the outcome of the analysis, is known as statistical influence analysis. Influence analysis examines the robustness of the model.

Influence analysis methodologies have been used extensively in classical linear models, and provided the basis for methodologies for use with LME models.

Computationally inexpensive diagnostics tools have been developed to examine the issue of influence (?). Studentized residuals, error contrast matrices and the inverse of the response variance covariance matrix are regular components of these tools.

15.2 Outline of Thesis

Thus the study of method comparison is introduced. The intention of this thesis is to progress the study of method comparison studies, using a statistical method known as Linear mixed effects models. Chapter two shall describe linear mixed effects models, and how the use of the linear mixed effects models have so far extended to method comparison studies. Implementations of important existing work shall be presented, using the R programming language.

Model diagnostics are an integral component of a complete statistical analysis. In chapter three model diagnostics shall be described in depth, with particular emphasis on linear mixed effects models, further to chapter two.

For the fourth chapter, important linear mixed effects model diagnostic methods shall be extended to method comparison studies, and proposed methods shall be demonstrated on data sets that have become well known in literature on method comparison. The purpose is to both calibrate these methods and to demonstrate applications for them. The last chapter shall focus on robust measures of important parameters such as agreement.

15.3 Extension of technique to LME Models

Model diagnostic techniques, well established for classical models, have since been adapted for use with linear mixed effects models. Diagnostic techniques for LME models are inevitably more difficult to implement, due to the increased complexity.

Beckman, Nachtsheim and Cook (1987) applied the local influence method of Cook (1986) to the analysis of the linear mixed model.

While the concept of influence analysis is straightforward, implementation in mixed models is more complex. Update formulae for fixed effects models are available only when the covariance parameters are assumed to be known.

If the global measure suggests that the points in U are influential, the nature of that influence should be determined. In particular, the points in U can affect

- the estimates of fixed effects

- the estimates of the precision of the fixed effects
- the estimates of the covariance parameters
- the estimates of the precision of the covariance parameters
- fitted and predicted values

15.3.1 Methods and Measures

The key to making deletion diagnostics useable is the development of efficient computational formulas, allowing one to obtain the case deletion diagnostics by making use of basic building blocks, computed only once for the full model.

Zewotir and Galpin (2005) lists several established methods of analyzing influence in LME models. These methods include

- Cook's distance for LME models,
- likelihood distance,
- the variance (information) ration,
- the Cook-Weisberg statistic,
- the Andrews-Prebigon statistic.

15.4 Haslett's Analysis

For fixed effect linear models with correlated error structure Haslett (1999) showed that the effects on the fixed effects estimate of deleting each observation in turn could be cheaply computed from the fixed effects model predicted residuals.

15.5 Augmented GLMs

Generalized linear models are a generalization of classical linear models.

The subscript M is a label referring to the mean model.

$$\begin{pmatrix} Y \\ \psi_M \end{pmatrix} = \begin{pmatrix} X & Z \\ 0 & I \end{pmatrix} \begin{pmatrix} \beta \\ \nu \end{pmatrix} + e^* \quad (15.1)$$

The error term e^* is normal with mean zero. The variance matrix of the error term is given by

$$\Sigma_a = \begin{pmatrix} \Sigma & 0 \\ 0 & D \end{pmatrix}. \quad (15.2)$$

$$X = \begin{pmatrix} T & Z \\ 0 & I \end{pmatrix} \delta = \begin{pmatrix} \beta \\ \nu \end{pmatrix} \quad (15.3)$$

$$y_a = T\delta + e^* \quad (15.4)$$

Weighted least squares equation

15.6 Covariance Parameters

The unknown variance elements are referred to as the covariance parameters and collected in the vector θ .

Chapter 16

Roy2013

<http://business.utsa.edu/wps/MSS/0017MSS-253-2013.pdf>

Testing the Equality of Mean Vectors for Paired Doubly Multivariate Observations

Example 2. (Mineral Data): This data set is taken from Johnson and Wichern (2007, p. 43). An investigator measured the mineral content of bones (radius, humerus and ulna) by photon absorptiometry to examine whether dietary supplements would slow bone loss in 25 older women. Measurements were recorded for three bones on the dominant and nondominant sides. Thus, the data is doubly multivariate and clearly $u = 2$ and $q = 3$. The bone mineral contents for the first 24 women one year after their participation in an experimental program is given in Johnson and Wichern (2007, p. 353).

Thus, for our analysis we take only first 24 women in the first data set. We test whether there has been a bone loss considering the data as doubly multivariate and has BCS structure. We rearrange the variables in the data set by grouping together the mineral content of the dominant sides of radius, humerus and ulna as the first three variables, that is, the variables in the first location ($u = 1$) and then the mineral contents for the non-dominant side of the same bones ($u = 2$)

16.1 Outlier Testing

A new outlier identification test for method comparison studies based on robust regression.

The identification of outliers in method comparison studies (MCS) is an important part of data analysis, as outliers can indicate serious errors in the measurement process. Common outlier tests

proposed in the literature usually require a homogeneous sample distribution and homoscedastic random error variances. However, datasets in MCS usually do not meet these assumptions. In this work, a new outlier test based on robust linear regression is proposed to overcome these special problems. The LORELIA (local reliability) residual test is based on a local, robust residual variance estimator, given as a weighted sum of the observed residuals. The new test is compared to a standard test proposed in the literature by a Monte Carlo simulation. Its performance is illustrated in examples.

16.2 Lorelia

Method comparison studies are performed in order to prove equivalence between two measurement methods or instruments. The identification of outliers is an important part of data analysis as outliers can indicate serious errors in the measurement process. Common outlier tests proposed in the literature require a homogeneous sample distribution and homoscedastic random error variances. However, datasets in method comparison studies usually do not meet these assumptions. To overcome this problem, different data transformation methods are proposed in the literature. However, they will only be applicable if the random errors can be described by simple additive or multiplicative models. In this work, a new outlier test based on robust linear regression is proposed which provides a general solution to the above problem. The LORELIA (LOcal RELIAbility) residual test is based on a local, robust residual variance estimator, given as a weighted sum of the observed residuals. Outlier limits are estimated from the actual data situation without making assumptions on the underlying error variance model. The performance of the new test is demonstrated in examples and simulations.

16.3 Note on Roy's paper

1. Basic model:

$$\begin{aligned} \mathbf{y}_i &= \mathbf{X}_i\boldsymbol{\beta} + \mathbf{Z}_i\mathbf{b}_i + \boldsymbol{\epsilon}_i, & i = 1, \dots, n \\ \mathbf{Z}_i &\sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}), & \boldsymbol{\epsilon}_i \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}) \end{aligned}$$

Assumptions are made about homoskedasticity.

2. General model:

$$\mathbf{y}_i = \mathbf{X}_i\boldsymbol{\beta} + \mathbf{Z}_i\mathbf{b}_i + \boldsymbol{\epsilon}_i, \quad i = 1, \dots, n$$

$$\mathbf{Z}_i \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Psi}), \quad \boldsymbol{\epsilon}_i \sim \mathcal{N}(\mathbf{0}, \sigma^2\boldsymbol{\Lambda})$$

Assumptions about homoskedasticity are relaxed (?, pg.202).

3. $\sigma^2\boldsymbol{\Lambda}$ is the general form for the VC structure for residuals.
4. The response vector \mathbf{y}_i comprises the observations of the subject, as measured by two methods, taking three measurements each. Hence a 6×1 random vector corresponding to the i th subject.

$$\mathbf{y}_i = (y_i^{j1}, y_i^{Jj2}, y_i^{j3}, y_i^{s1}, y_i^{s2}, y_i^{s3})' \quad (16.1)$$

5. The number of replicates is p . A subject will have up to $2p$ measurements, for the two instrument case, i.e. $\text{Max}(n_i) = 2p$. (Let k denote number of instruments, which is assumed to be 2 unless stated otherwise.) For the blood pressure data $p = 3$.

16.3.1 Outliers and Leverage

The question of whether or not a point should be considered an outlier must also be addressed. An outlier is an observation whose true value is unusual given its value on the predictor variables. The leverage of an observation is a further consideration. Leverage describes an observation with an extreme value on a predictor variable is a point with high leverage. High leverage points can have a great amount of effect on the estimate of regression coefficients.

Influence can be thought of as the product of leverage and outlierness. An observation is said to be influential if removing the observation substantially changes the estimate of the regression coefficients. The R programming language has a variety of methods used to study each of the aspects for a linear model. While linear models and GLMS can be studied with a wide range of well-established diagnostic techniques, the choice of methodology is much more restricted for the case of LMEs.

For classical linear models, residual diagnostics are typically conducted using a plot of the observed residuals and the predicted values. A visual inspection for the presence of trends inform the analyst on the validity of distributional assumptions, and to detect outliers and influential observations.

16.3.2 Matrix Notation for Case Deletion

For notational simplicity, $\mathbf{A}(i)$ denotes an $n \times m$ matrix \mathbf{A} with the i -th row removed, a_i denotes the i -th row of \mathbf{A} , and a_{ij} denotes the (i, j) -th element of \mathbf{A} .

16.3.3 Extension of Diagnostic Methods to LME models

When similar notions of statistical influence are applied to mixed models, things are more complicated. Removing data points affects fixed effects and covariance parameter estimates. Update formulas for *leave-one-out* estimates typically fail to account for changes in covariance parameters.

? noted the case deletion diagnostics techniques have not been applied to linear mixed effects models and seeks to develop methodologies in that respect. ? develops these techniques in the context of REML.

? noted the case deletion diagnostics techniques had not been applied to linear mixed effects models and seeks to develop methodologies in that respect. ? develops these techniques in the context of REML. [~~~~~ origin/master](#)

Demidenko (2004) extends several regression diagnostic techniques commonly used in linear regression, such as leverage, infinitesimal influence, case deletion diagnostics, Cook's distance, and local influence to the linear mixed-effects model. In each case, the proposed new measure has a direct interpretation in terms of the effects on a parameter of interest, and reduces to the familiar linear regression measure when there are no random effects.

The new measures that are proposed by Demidenko (2004) are explicitly defined functions and do not require re-estimation of the model, especially for cluster deletion diagnostics. The basis for both the cluster deletion diagnostics and Cook's distance is a generalization of Miller's simple update formula for case deletion for linear models. Furthermore Demidenko (2004) shows how Pregibon's infinitesimal case deletion diagnostics is adapted to the linear mixed-effects model.

Demidenko (2004) proposes two likelihood-based diagnostics for identifying individuals that can influence the choice between two competing models.

16.4 Model Validation Framework

In statistical modelling, the process of model validation is a critical step of model fitting process, but also a step that is too often overlooked. A very simple procedure is to examine commonly-used metrics, such as the R^2 value. However, using a small handful of simple measures and methods is insufficient to properly assess the quality of a fitted model. To do so properly, a full and comprehensive analysis that tests of all of the assumptions, as far as possible, must be carried out.

? describes the model validation framework as comprised of the following tasks

- overall measures of goodness-of-fit
- the informal, graphical examination of estimates of model errors to assess the quality of distributional assumptions: residual analysis
- the quantitative assessment of the inter-relationship of model components; for example, collinearity diagnostics
- the qualitative and quantitative assessment of influence of cases on the analysis, i.e. influence analysis.

The sensitivity of a model is studied through measures that express its stability under perturbations. You are not interested in a model that is either overly stable or overly sensitive. Changes in the data or model components should produce commensurate changes in the model output. The difficulty is to determine when the changes are substantive enough to warrant further investigation, possibly leading to a reformulation of the model or changes in the data (such as dropping outliers). This paper is primarily concerned with stability of linear mixed models to perturbations of the data; that is, with influence analysis.

16.4.1 Residual Analysis

Residual analysis is a widely used model validation technique. A residual is simply the difference between an observed value and the corresponding fitted value, as predicted by the model. The rationale is that, if the model is properly fitted to the model, then the residuals would approximate the random errors that one should expect. that is to say, if the residuals behave randomly, with no discernible trend, the model has fitted the data well. If some sort of non-random trend is evident in the model, then the model can be considered to be poorly fitted.

Statistical software environments, such as the R Programming language, provides a suite of tests and graphical procedure sfor appraising a fitted linear model, with several of these procedures analysing the model residuals.

16.4.2 Outliers and Leverage

The question of whether or not a point should be considered an outlier must also be addressed. An outlier is an observation whose true value is unusual given its value on the predictor variables. The leverage of an observation is a further consideration. Leverage describes an observation with an extreme value on a predictor variable is a point with high leverage. High leverage points can have a great amount of effect on the estimate of regression coefficients.

Influence can be thought of as the product of leverage and outlierness. An observation is said to be influential if removing the observation substantially changes the estimate of the regression coefficients. The R programming language has a variety of methods used to study each of the aspects for a linear model. While linear models and GLMS can be studied with a wide range of well-established diagnostic techniques, the choice of methodology is much more restricted for the case of LMEs.

16.5 Regression Of Differences On Averages

Further to Carstensen, we can formulate the two measurements y_1 and y_2 as follows:

$$y_1 = \alpha + \beta\mu + \epsilon_1$$

$$y_2 = \alpha + \beta\mu + \epsilon_2$$

16.5.1 Note 1: Coefficient of Repeatability

The coefficient of repeatability is a measure of how well a measurement method agrees with itself over replicate measurements (Bland and Altman, 1999). Once the within-item variability is known, the computation of the coefficients of repeatability for both methods is straightforward.

16.5.2 Note 2: Carstensen model in the single measurement case

Carstensen (2004) presents a model to describe the relationship between a value of measurement and its real value. The non-replicate case is considered first, as it is the context of the Bland-Altman plots. This model assumes that inter-method bias is the only difference between the two methods.

$$y_{mi} = \alpha_m + \mu_i + e_{mi} \quad e_{mi} \sim \mathcal{N}(0, \sigma_m^2) \quad (16.2)$$

The differences are expressed as $d_i = y_{1i} - y_{2i}$.

For the replicate case, an interaction term c is added to the model, with an associated variance component.

16.5.3 Note 3: Model terms

It is important to note the following characteristics of this model.

- Let the number of replicate measurements on each item i for both methods be n_i , hence $2 \times n_i$ responses. However, it is assumed that there may be a different number of replicates made for different items. Let the maximum number of replicates be p . An item will have up to $2p$ measurements, i.e. $\max(n_i) = 2p$.
- Later on \mathbf{X}_i will be reduced to a 2×1 matrix, to allow estimation of terms. This is due to a shortage of rank. The fixed effects vector can be modified accordingly.
- \mathbf{Z}_i is the $2n_i \times 2$ model matrix for the random effects for measurement methods on item i .
- \mathbf{b}_i is the 2×1 vector of random-effect coefficients on item i , one for each method.

- ϵ is the $2n_i \times 1$ vector of residuals for measurements on item i .
- \mathbf{G} is the 2×2 covariance matrix for the random effects.
- \mathbf{R}_i is the $2n_i \times 2n_i$ covariance matrix for the residuals on item i .
- The expected value is given as $E(\mathbf{y}_i) = \mathbf{X}_i\boldsymbol{\beta}$. (Hamlett et al., 2004)
- The variance of the response vector is given by $\text{Var}(\mathbf{y}_i) = \mathbf{Z}_i\mathbf{G}\mathbf{Z}_i' + \mathbf{R}_i$ (Hamlett et al., 2004).

Roy's uses and LME model approach to provide a set of formal tests for method comparison studies.

Four candidate models are fitted to the data.

These models are similar to one another, but for the imposition of equality constraints.

These tests are the pairwise comparison of candidate models, one formulated without constraints, the other with a constraint.

Roy's model uses fixed effects $\beta_0 + \beta_1$ and $\beta_0 + \beta_1$ to specify the mean of all observations by methods 1 and 2 respectively.

Roy adheres to Random Effect ideas in ANOVA

Roy treats items as a sample from a population.

Allocation of fixed effects and random effects are very different in each model

Carstensen's interest lies in the difference between the population from which they were drawn.

Carstensen's model is a mixed effects ANOVA.

$$Y_{mir} = \alpha_m + \mu_i + c_{mi} + e_{mir}, \quad c_{mi} \sim \tau_{\updownarrow}^{\epsilon}, \quad e_{mir} \sim \sigma_{\updownarrow}^{\epsilon},$$

This model includes a method by item interaction term.

Carstensen presents two models. One for the case where the replicates, and a second for when they are linked.

Carstensen's model does not take into account either between-item or within-item covariance between methods.

In the presented example, it is shown that Roy's LoAs are lower than those of Carstensen. Carstensen makes some interesting remarks in this regard.

The only slightly non-standard (meaning "not often used") feature is the differing residual variances between methods.

16.5.4 Note 3: Model terms

It is important to note the following characteristics of this model.

- Let the number of replicate measurements on each item i for both methods be n_i , hence $2 \times n_i$ responses. However, it is assumed that there may be a different number of replicates made for different items. Let the maximum number of replicates be p . An item will have up to $2p$ measurements, i.e. $\max(n_i) = 2p$.
- Later on \mathbf{X}_i will be reduced to a 2×1 matrix, to allow estimation of terms. This is due to a shortage of rank. The fixed effects vector can be modified accordingly.
- \mathbf{Z}_i is the $2n_i \times 2$ model matrix for the random effects for measurement methods on item i .
- \mathbf{b}_i is the 2×1 vector of random-effect coefficients on item i , one for each method.
- $\boldsymbol{\epsilon}$ is the $2n_i \times 1$ vector of residuals for measurements on item i .
- \mathbf{G} is the 2×2 covariance matrix for the random effects.
- \mathbf{R}_i is the $2n_i \times 2n_i$ covariance matrix for the residuals on item i .
- The expected value is given as $E(\mathbf{y}_i) = \mathbf{X}_i\boldsymbol{\beta}$. (Hamlett et al., 2004)
- The variance of the response vector is given by $\text{Var}(\mathbf{y}_i) = \mathbf{Z}_i\mathbf{G}\mathbf{Z}_i' + \mathbf{R}_i$ (Hamlett et al., 2004).

16.6 Regression Of Differences On Averages

Further to Carstensen, we can formulate the two measurements y_1 and y_2 as follows:

$$y_1 = \alpha + \beta\mu + \epsilon_1$$

$$y_2 = \alpha + \beta\mu + \epsilon_2$$

For classical linear models, residual diagnostics are typically conducted using a plot of the observed residuals and the predicted values. A visual inspection for the presence of trends inform the analyst on the validity of distributional assumptions, and to detect outliers and influential observations.

Appendix to Section 4

As an appendix to section 4, an appraisal of the current state of development (or lack thereof) for current implemenations for LME models, particularly for `nlme` and `lme4` fitted models.

Crucially, a review of internet resources indicates that almost all of the progress in this regard has been done for `lme4` fitted models, specifically the *Influence.ME* R package. (Nieuwenhuis et 2012)

Conversely there is very little for `nlme` models. To delve into this mor, one would immediately investigate the current development workflow for both packages.

As an aside, Douglas Bates was arguably the most prominent R developer working in the LME area. However Bates has now prioritised the development of LME models in another computing environment , i.e Julia.

Important Consideration for MCS

The key issue is that `nlme` allows for the particular specification of Roy’s Model, speciifcally direct spefification of the VC matrices for within subject and between subject residuals. The `lme4` package does not allow for this. To advance the ideas that eminate from Roys’ paper, one is required to use the `nlme` context. However, to take advantage of the infrastructure already provided for `lme4` models, one may change the research question away from that of Roy’s paper. To this end, an exploration of what textitinfluence.ME can accomplished is merited. As a complement to this, one can also consider how to properly employ the R^2 measure, in the context of Methoc Comparison Studies, further to the work by Edwards et al, namely “An R^2 statistic for fixed effects in the linear mixed model”.

Abstract for “An R^2 statistic for fixed effects in the linear mixed model”

Statisticians most often use the linear mixed model to analyze Gaussian longitudinal data.

The value and familiarity of the R^2 statistic in the linear univariate model naturally creates great interest in extending it to the linear mixed model. We define and describe how to compute a model R^2 statistic for the linear mixed model by using only a single model.

The proposed R^2 statistic measures multivariate association between the repeated outcomes and the fixed effects in the linear mixed model. The R^2 statistic arises as a 11 function of an appropriate F statistic for testing all fixed effects (except typically the intercept) in a full model.

The statistic compares the full model with a null model with all fixed effects deleted (except typically the intercept) while retaining exactly the same covariance structure.

Furthermore, the R^2 statistic leads immediately to a natural definition of a partial R^2 statistic. A mixed model in which ethnicity gives a very small p-value as a longitudinal predictor of blood pressure (BP) compellingly illustrates the value of the statistic.

In sharp contrast to the extreme p-value, a very small R^2 , a measure of statistical and scientific importance, indicates that ethnicity has an almost negligible association with the repeated BP outcomes for the study.

The nlme package

With regards to `nlme`, the torch has been passed to Galecki Galecki & Burzykowski (UMich. and Hasselt respecitely). Galecki & Burzykowski published *Linear Mixed Effects Models using R*. Also, the accompanying R package, `nlmeU` package is under current development, with a version being released XXXX.

The lme4 package

The `lme4` package is also under active development, under the leadership of Ben Bolker (McMaster University). According to CRAN, the LME4 package, fits linear and generalized linear mixed-effects models

The models and their components are represented using S4 classes and methods. The core computational algorithms are implemented using the Eigen C++ library for numerical linear algebra and RcppEigen "glue". (CRAN)

The key issue is that `nlme` allows for the particular specification of Roy's Model, specifically direct specification of the VC matrices for within subject and between subject residuals. The `lme4` package does not allow for this. To advance the ideas that emanate from Roys' paper, one is required to use the `nlme` context. However, to take advantage of the infrastructure already provided for `lme4` models, one may change the research question away from that of Roy's paper. To this end, an exploration of what textitinfluence.ME can accomplished is merited. As a complement to this, one can also consider how to properly employ the R^2 measure, in the context of Methoc Comparison Studies, further to the work by Edwards et al, namely "An R^2 statistic for fixed effects in the linear mixed model".

Abstract for “An R^2 statistic for fixed effects in the linear mixed model”

Statisticians most often use the linear mixed model to analyze Gaussian longitudinal data.

The value and familiarity of the R^2 statistic in the linear univariate model naturally creates great interest in extending it to the linear mixed model. We define and describe how to compute a model R^2 statistic for the linear mixed model by using only a single model.

The proposed R^2 statistic measures multivariate association between the repeated outcomes and the fixed effects in the linear mixed model. The R^2 statistic arises as a 11 function of an appropriate F statistic for testing all fixed effects (except typically the intercept) in a full model.

The statistic compares the full model with a null model with all fixed effects deleted (except typically the intercept) while retaining exactly the same covariance structure.

Furthermore, the R^2 statistic leads immediately to a natural definition of a partial R^2 statistic. A mixed model in which ethnicity gives a very small p-value as a longitudinal predictor of blood pressure (BP) compellingly illustrates the value of the statistic.

In sharp contrast to the extreme p-value, a very small R^2 , a measure of statistical and scientific importance, indicates that ethnicity has an almost negligible association with the repeated BP outcomes for the study.

$$r_{mi} = x_i^T \hat{\beta} \tag{16.3}$$

16.6.1 Marginal Residuals

$$\begin{aligned}\hat{\beta} &= (X^T R^{-1} X)^{-1} X^T R^{-1} Y \\ &= BY\end{aligned}$$

16.7 Covariance Parameters

The unknown variance elements are referred to as the covariance parameters and collected in the vector θ .

16.7.1 Methods and Measures

The key to making deletion diagnostics useable is the development of efficient computational formulas, allowing one to obtain the case deletion diagnostics by making use of basic building blocks, computed only once for the full model.

Zewotir and Galpin (2005) lists several established methods of analyzing influence in LME models. These methods include

- Cook's distance for LME models,
- likelihood distance,
- the variance (information) ration,
- the Cook-Weisberg statistic,
- the Andrews-Prebigon statistic.

16.8 Missing Data in Method Comparison Studies

The matter of missing data has not been commonly encountered in either Method Comparison Studies or Linear Mixed Effects Modelling. However Roy (2009) deals with the relevant assumptions regarding missing data.

Galecki & Burzykowski (2013) tackles the subject of missing data in LME Modelling.

Furthermore the nlmeU package includes the `patMiss` function, which “allows to compactly present pattern of missing data in a given vector/matrix/data frame or combination of thereof”.

16.9 Leave-One-Out Diagnostics with lmeU

Galecki et al discuss the matter of LME influence diagnostics in their book, although not into great detail.

The command `lmeU` fits a model with a particular subject removed. The identifier of the subject to be removed is passed as the only argument

A plot of the per-observation diagnostics individual subject log-likelihood contributions can be rendered.

- R command and R object - Typewriter Font
- R Package name - Italics
- Selected Acronyms and Proper Nouns - Italics

- This chapter is broken into two parts. The first part is a review of diagnostics methods for linear models, intended to acquaint the reader with the subject, and also to provide a basis for material covered in the second part. Particular attention is drawn to graphical methods.
- The second part of the chapter looks at diagnostics techniques for LME models, firstly covering the theory, then proceeding to a discussion on implementing these using R code.
- While a substantial body of work has been developed in this area, there are still areas worth exploring. In particular the development of graphical techniques pertinent to LME models should be looked at.

16.10 Introduction

Chapter 17

Model Diagnostics

Abstract

This chapter is broken into two parts. The first part is a review of diagnostics methods for linear models, intended to acquaint the reader with the subject, and also to provide a basis for material covered in the second part. Particular attention is drawn to graphical methods.

17.1 Carstensen

Let y_{mir} be the r th replicate measurement on the i th item by the m th method, where $m = 1, 2$, $i = 1, \dots, N$, and $r = 1, \dots, n_i$. When the design is balanced and there is no ambiguity we can set $n_i = n$. The LME model underpinning Roy's approach can be written

$$y_{mir} = \beta_0 + \beta_m + b_{mi} + \epsilon_{mir}. \quad (17.1)$$

Here β_0 and β_m are fixed-effect terms representing, respectively, a model intercept and an overall effect for method m . The β terms can be gathered together into (fixed effect) intercept terms $\alpha_m = \beta_0 + \beta_m$. The b_{1i} and b_{2i} terms are correlated random effect parameters having $E(b_{mi}) = 0$ with $\text{Var}(b_{mi}) = g_m^2$ and $\text{Cov}(b_{mi}, b_{m'i}) = g_{12}$. The random error term for each response is denoted ϵ_{mir} having $E(\epsilon_{mir}) = 0$, $\text{Var}(\epsilon_{mir}) = \sigma_m^2$, $\text{Cov}(b_{mir}, b_{m'ir}) = \sigma_{12}$, $\text{Cov}(\epsilon_{mir}, \epsilon_{mir'}) = 0$ and $\text{Cov}(\epsilon_{mir}, \epsilon_{m'ir'}) = 0$. Two methods of measurement are in complete agreement if the null hypotheses $H_1: \beta_1 = \beta_2$ and $H_2: \sigma_1^2 = \sigma_2^2$ and

$H_3: g_1^2 = g_2^2$ hold simultaneously. ? proposes a Bonferroni correction to control the familywise error rate for tests of $\{H_1, H_2, H_3\}$ and account for difficulties arising due to multiple testing. Let $\omega_m^2 = \sigma_m^2 + g_m^2$ represent the overall variability of method m . Roy also integrates H_2 and H_3 into a single testable hypothesis $H_4: \omega_1^2 = \omega_2^2$. CONCERNS?

? demonstrates how to implement a method comparison study further to model (1) using the SAS proc mixed package. Carstensen et al. (2008) demonstrates how to construct limits of agreement using SAS, STATA and R. In the case of SAS, the PROC MIXED procedure is used. Implementation in R is performed using the nlme package (?).

Carstensen et al. (2008) remarks that the implementation using R is quite “arcane”.

As R is freely available, this paper demonstrates an implementation of Roy’s model using R.

The R statistical software package is freely available.

The LME model is very easy to implement using PROC MIXED of SAS and the results are also easy to interpret. The SAS proc mixed procedure has very simple syntax.

As the required code to fit the models is complex, R code necessary to fit the models is provided.

A demonstration is provided on how to use the output to perform the tests, and to compute limits of agreement.

We assume the data are formatted as a dataset with four columns named:

meth, method of measurement, the number of methods being M, item, items (persons, samples) measured by each method, of which there are I, repl, replicate indicating repeated measurement of the same item by the same method, and y, the measurement.

17.1.1 Remarks on the Multivariate Normal Distribution

Diligence is required when considering the models. Carstensen specifies his models in terms of the univariate normal distribution. Roy’s model is specified using the bivariate normal distribution. This gives rises to a key difference between the two model, in that a bivariate model accounts for covariance between the variables of interest. The multivariate normal distribution of a k -dimensional random vector $X = [X_1, X_2, \dots, X_k]$ can be written in the following notation:

$$X \sim \mathcal{N}(\mu, \Sigma),$$

or to make it explicitly known that X is k -dimensional,

$$X \sim \mathcal{N}_k(\mu, \Sigma).$$

with k -dimensional mean vector

$$\mu = [E[X_1], E[X_2], \dots, E[X_k]]$$

and $k \times k$ covariance matrix

$$\Sigma = [\text{Cov}[X_i, X_j]], \quad i = 1, 2, \dots, k; \quad j = 1, 2, \dots, k$$

1. Univariate Normal Distribution

$$X \sim \mathcal{N}(\mu, \sigma^2),$$

2. Bivariate Normal Distribution

(a)

$$X \sim \mathcal{N}_2(\mu, \Sigma),$$

(b)

$$\mu = \begin{pmatrix} \mu_x \\ \mu_y \end{pmatrix}, \quad \Sigma = \begin{pmatrix} \sigma_x^2 & \rho\sigma_x\sigma_y \\ \rho\sigma_x\sigma_y & \sigma_y^2 \end{pmatrix}.$$

17.2 Modelling Agreement with LME Models

Roys uses an LME model approach to provide a set of formal tests for method comparison studies.

Four candidate models are fitted to the data.

These models are similar to one another, but for the imposition of equality constraints.

These tests are the pairwise comparison of candidate models, one formulated without constraints, the other with a constraint.

Roy's model uses fixed effects $\beta_0 + \beta_1$ and $\beta_0 + \beta_1$ to specify the mean of all observations by methods 1 and 2 respectively.

Roy adheres to Random Effect ideas in ANOVA

Roy treats items as a sample from a population.

Allocation of fixed effects and random effects are very different in each model

Carstensen's interest lies in the difference between the population from which they were drawn.

Carstensen's model is a mixed effects ANOVA.

$$Y_{mir} = \alpha_m + \mu_i + c_{mi} + e_{mir}, \quad c_{mi} \sim \tau_{\Downarrow}^{\epsilon}, \quad e_{mir} \sim \sigma_{\Downarrow}^{\epsilon},$$

This model includes a method by item iteration term.

Carstensen presents two models. One for the case where the replicates, and a second for when they are linked.

Carstensen's model does not take into account either between-item or within-item covariance between methods.

In the presented example, it is shown that Roy's LoAs are lower than those of Carstensen. Carstensen makes some interesting remarks in this regard.

The only slightly non-standard (meaning "not often used") feature is the differing residual variances between methods.

It is also desirable to measure the influence of the case deletions on the covariance matrix of $\hat{\beta}$.

- *The previous Section (Section 4) is a literary review of residual diagnostics and influence procedures for Linear Mixed Effects Models, drawing heavily on Schabenberger and Zewotir.*
- *Section 4 begins with an introduction to key topics in residual diagnostics, such as influence, leverage, outliers and Cook's distance. Other concepts such as DFFITS and DFBETAs will be introduced briefly, mostly to explain why they are not particularly useful for the Method Comparison context, and therefore are not elaborated upon.*
- *In brief, Variable Selection is not applicable to Method Comparison Studies, in the commonly used context. Testing a rather simplistic specified model against one with more random effects terms is tractable, but this research question is of secondary importance.*

17.2.1 Matrix Notation for Case Deletion

For notational simplicity, $\mathbf{A}(i)$ denotes an $n \times m$ matrix \mathbf{A} with the i -th row removed, a_i denotes the i -th row of \mathbf{A} , and a_{ij} denotes the (i, j) -th element of \mathbf{A} .

17.3 Extension of Diagnostic Methods to LME models

When similar notions of statistical influence are applied to mixed models, things are more complicated. Removing data points affects fixed effects and covariance parameter estimates. Update formulas for *leave-one-out* estimates typically fail to account for changes in covariance parameters.

? noted the case deletion diagnostics techniques have not been applied to linear mixed effects models and seeks to develop methodologies in that respect. ? develops these techniques in the context of REML.

? noted the case deletion diagnostics techniques had not been applied to linear mixed effects models and seeks to develop methodologies in that respect. ? develops these techniques in the context of REML.

Demidenko (2004) extends several regression diagnostic techniques commonly used in linear regression, such as leverage, infinitesimal influence, case deletion diagnostics, Cook's distance, and local

influence to the linear mixed-effects model. In each case, the proposed new measure has a direct interpretation in terms of the effects on a parameter of interest, and reduces to the familiar linear regression measure when there are no random effects.

The new measures that are proposed by Demidenko (2004) are explicitly defined functions and do not require re-estimation of the model, especially for cluster deletion diagnostics. The basis for both the cluster deletion diagnostics and Cook's distance is a generalization of Miller's simple update formula for case deletion for linear models. Furthermore Demidenko (2004) shows how Pregibon's infinitesimal case deletion diagnostics is adapted to the linear mixed-effects model.

Demidenko (2004) proposes two likelihood-based diagnostics for identifying individuals that can influence the choice between two competing models.

17.3.1 Remarks on the Multivariate Normal Distribution

Diligence is required when considering the models. Carstensen specifies his models in terms of the univariate normal distribution. Roy's model is specified using the bivariate normal distribution. This gives rises to a key difference between the two model, in that a bivariate model accounts for covariance between the variables of interest. The multivariate normal distribution of a k -dimensional random vector $X = [X_1, X_2, \dots, X_k]$ can be written in the following notation:

$$X \sim \mathcal{N}(\mu, \Sigma),$$

or to make it explicitly known that X is k -dimensional,

$$X \sim \mathcal{N}_k(\mu, \Sigma).$$

with k -dimensional mean vector

$$\mu = [E[X_1], E[X_2], \dots, E[X_k]]$$

and $k \times k$ covariance matrix

$$\Sigma = [\text{Cov}[X_i, X_j]], \ i = 1, 2, \dots, k; \ j = 1, 2, \dots, k$$

1. Univariate Normal Distribution

$$X \sim \mathcal{N}(\mu, \sigma^2),$$

2. Bivariate Normal Distribution

(a)

$$X \sim \mathcal{N}_2(\mu, \Sigma),$$

(b)

$$\mu = \begin{pmatrix} \mu_x \\ \mu_y \end{pmatrix}, \quad \Sigma = \begin{pmatrix} \sigma_x^2 & \rho\sigma_x\sigma_y \\ \rho\sigma_x\sigma_y & \sigma_y^2 \end{pmatrix}.$$

17.4 Repeated measurements in LME models

In many statistical analyzes, the need to determine parameter estimates where multiple measurements are available on each of a set of variables often arises. Further to Lam et al. (1999), Hamlett et al. (2004) performs an analysis of the correlation of replicate measurements, for two variables of interest, using LME models.

Let y_{Aij} and y_{Bij} be the j th repeated observations of the variables of interest A and B taken on the i th subject. The number of repeated measurements for each variable may differ for each individual. Both variables are measured on each time points. Let n_i be the number of observations for each variable, hence $2 \times n_i$ observations in total.

It is assumed that the pair y_{Aij} and y_{Bij} follow a bivariate normal distribution.

$$\begin{pmatrix} y_{Aij} \\ y_{Bij} \end{pmatrix} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}) \text{ where } \boldsymbol{\mu} = \begin{pmatrix} \mu_A \\ \mu_B \end{pmatrix}$$

The matrix $\boldsymbol{\Sigma}$ represents the variance component matrix between response variables at a given time point j .

$$\boldsymbol{\Sigma} = \begin{pmatrix} \sigma_A^2 & \sigma_{AB} \\ \sigma_{AB} & \sigma_B^2 \end{pmatrix}$$

σ_A^2 is the variance of variable A , σ_B^2 is the variance of variable B and σ_{AB} is the covariance of the two variable. It is assumed that $\boldsymbol{\Sigma}$ does not depend on a particular time point, and is the same over all time points.

Bibliography

- ACR (2008). Acute Chest Pain (suspected aortic dissection) - American College of Radiology Expert Group Report.
- Akaike, H. (1974). A new look at the statistical model identification. *IEEE Transactions on Automatic Control* 19(6), 716–723.
- Altman, D. and J. Bland (1983). Measurement in medicine: The analysis of method comparison studies. *Journal of the Royal Statistical Society. Series D (The Statistician)* 32(3), 307–317.
- Barnhart, H., M. Haber, and L. Lin (2007). An overview of assessing agreement with continuous measurements. *Journal of Biopharmaceutical Statistics* 17, 529–569.
- Bartko, J. (1994). Measures of agreement: A single procedure. *Statistics in Medicine* 13, 737–745.
- Benjamini, Y. and Y. Hochberg (1995). Controlling the false discovery rate: A practical and powerful approach to multiple testing. *J. Roy. Statist. Soc. Ser. B* 57(1), 289–300.
- Bland, J. and D. Altman (1986). Statistical methods for assessing agreement between two methods of clinical measurement. *The Lancet* i, 307–310.
- Bland, J. and D. Altman (1995). Comparing methods of measurement - why plotting difference against standard method is misleading. *The Lancet* 346, 1085–87.
- Bland, J. and D. Altman (1999). Measuring agreement in method comparison studies. *Statistical Methods in Medical Research* 8(2), 135–160.
- Brown, H. and R. Prescott (1999). *Applied Mixed Models In Medicine*. John Wiley and Sons.

- Carstensen, B. (2004). Comparing and predicting between several methods of measurement. *Biostatistics* 5(3), 399–413.
- Carstensen, B. (2010). *Comparing Clinical Measurement Methods : A Practical Guide*. Wiley.
- Carstensen, B., J. Simpson, and L. C. Gurrin (2008). Statistical models for assessing agreement in method comparison studies with replicate measurements. *The International Journal of Biostatistics* 4(1).
- Cook, R. (1977). Detection of influential observations in linear regression. *Technometrics* 19, 15–18.
- Cornbleet, P. J. and D. Cochrane (1979). Regression methods for assessing agreement between two methods of clinical measurement. *Journal of Clinical Chemistry* 24(2), 342–345.
- Demidenko, E. (2004). *Mixed Models: Theory And Application*. Dartmouth College: Wiley Interscience.
- Dempster, A., N. Laird, and D. Rubin (1977). Maximum likelihood from incomplete data via the em algorithm. *ournal of the Royal Statistical Society. Series B* 39(1), 1–38.
- Dewitte, K., C. Fierens, D. Stckl, and L. M. Thienpont (2002). Application of the Bland Altman plot for interpretation of method-comparison studies: A critical investigation of its practice. *Clinical Chemistry* 48, 799–801.
- Dodge, Y. (2003). *The Oxford Dictionary of Statistical Terms*. Oxford University Press.
- Dunn, G. (2002). *Statistical Evaluation of Measurement Error* (Second ed.). Stanford: American Mathematical Society and Digital Press.
- Eisenhart, C. (1947). The assumptions underlying the analysis of variance. *Biometrics* 3(1), 1–21.
- Fisher, R. (1918). The correlation between relatives on the supposition of mendelian inheritance. *Transactions of the Royal Society of Edinburgh* 2, 399–433.
- Grubbs, F. (1973). Errors of measurement, precision, accuracy and the statistical comparison of measuring instruments. *Technometrics* 15(1), 53–66.

- Hamlett, A., L. Ryan, and R. Wolfinger (2004). On the use of PROC MIXED to estimate correlation in the presence of repeated measures. *Proceedings of the Statistics and Data Analysis Section, SAS Users Group International 198-229*, 1–7.
- Hartley, H. and J. Rao (1967). Maximum-likelihood estimation for the mixed analysis of variance model. *Biometrika* 54(1/2), 93–108.
- Harville, D. A. (1977). Maximum likelihood approaches to variance component estimation and to related problems. *Journal of American Statistical Association* 72(358), 320–338.
- Haslett, J. (1999). A simple derivation of deletion diagnostic results for the general linear model with correlated errors. *Journal of the Royal Statistical Society (Series B)* 61, 603–609.
- Haslett, J. and D. Dillane (2004). Application of ‘delete = replace’ to deletion diagnostics for variance component estimation. *Journal of the Royal Statistical Society (Series B)* 66, 131–143.
- Haslett, J. and K. Hayes (1998). Residuals for the linear model with general covariance structure. *Journal of the Royal Statistical Society (Series B)* 60, 201–215.
- Henderson, C. (1953). Estimation of variance and covariance components. *Biometrics* 9(2), 226–252.
- Henderson, C., O. Kempthorne, S. Searle, and C. von Krosigk (1959). The estimation of environmental and genetic trends from records subject to culling. *Biometrics* 15, 192–218.
- Henderson, C., O. Kempthorne, S. Searle, and C. von Krosigk (1963). The estimation of environmental and genetic trends from records subject to culling. *Biometrics* 15, 192–218.
- Henderson, C., O. Kempthorne, S. Searle, and C. von Krosigk (1973). The estimation of environmental and genetic trends from records subject to culling. *Biometrics* 15, 192–218.
- Henderson, C., O. Kempthorne, S. Searle, and C. von Krosigk (1984). The estimation of environmental and genetic trends from records subject to culling. *Biometrics* 15, 192–218.
- Henderson, C. R. (1950). Estimation of genetic parameters (abstract). *Annals of Mathematical Statistics* 21, 309–310.

- Laird, N. M. and J. H. Ware (1982). Random-effects models for longitudinal data. *Biometrics* 38(4), 963–974.
- Lam, M., K. Webb, and D. O'Donnell (1999). Correlation between two variables in repeated measurements. *American Statistical Association, Proceedings of the Biometric Session*, 213–218.
- Lewis, P., P. Jones, J. Polak, and H. Tillitson (1991). The problem of conversion in method comparison studies. *Applied Statistics* 40, 105–112.
- Lin, L. (2000). Total deviation index for measuring individual agreement with applications in laboratory performance and bioequivalence. *Statistics in medicine* 97, 255–270.
- Lin, L., A. Hedayat, B. Sinha, and M. Yang (2002). Statistical methods in assessing agreement: Models issues and tools. *Journal of American Statistical Association* 97, 257–270.
- Lin, S. C., D. M. Whipple, and Charles S Ho (1991). Evaluation of statistical equivalence using limits of agreement and associates sample size calculation. *Communications in Statistics - Theory and Methods* 27(6), 1419–1432.
- Linnet, K. (1998). Performance of deming regression analysis in case of misspecified analytical error ratio in method comparison studies. *Clinical Chemistry* 44, 1024–1031.
- Ludbrook, J. (1997). Comparing methods of measurement. *Clinical and Experimental Pharmacology and Physiology* 24, 193–203.
- Ludbrook, J. (2002). Statistical techniques for comparing measurers and methods of measurement: a critical review. *Clinical and Experimental Pharmacology and Physiology* 29, 527–536.
- Luiz, R., A. Costa, P. Kale, and G. Werneck (2003). Assessment of agreement of a quantitative variable: a new graphical approach. *Journal of Clinical Epidemiology* 56, 963–967.
- Mantha, S., M. F. Roizen, L. A. Fleisher, R. Thisted, and J. Foss (2000). Comparing methods of clinical measurement: Reporting standards for bland and altman analysis. *Anaesthesia and Analgesia* 90, 593–602.

- McCullough, C. and S. Searle (2001). *Generalized , Linear and Mixed Models*. Wiley Interscience.
- Myles, P. (2007). Using the Bland Altman method to measure agreement with repeated measures. *British Journal of Anaesthesia* 99(3), 309–311.
- NIST (2009). Cesium fountain atomic clock : The primary time and frequency standard for the United States. <http://tf.nist.gov/timefreq/cesium/fountain.htm>.
- Paterson, H. and R. Thompson (1971). Recovery of inter-block information when block sizes are unequal. *Biometrika* 58(3), 545–554.
- Patterson, H. and R. Thompson (1971). Recovery of inter-block information when block sizes are unequal. *Biometrika* 58(3), 545–554.
- Phelps, C. and A. Hutson (1995). Estimating diagnostic test accuracy using a fuzzy gold standard. *Medical decision making* 15, 144–57.
- Pinheiro, J. and D. Bates (1994). *Mixed Effects Models in S and S plus* (2nd ed.). Reading, Massachusetts: Springer.
- Pizzi, N. (1999). Fuzzy pre-processing of gold standards as applied to biomedical spectra classification. *Artificial Intelligence in Medicine* 16, 171–182.
- Roy, A. (2009). An application of the linear mixed effects model to ass the agreement between two methods with replicated observations. *Journal of Biopharmaceutical Statistics* 19, 150–173.
- Schabenberger, O. (2004). Mixed model influence diagnostics. 18929.
- Searle, S. (1997). *Linear Models*. Wiley classics Library.
- Singer, J. D. (1998). Using sas proc mixed to fit multilevel models, hierarchical models, and individual growth models. *Journal of Educational and Behavioral Statistics* 24(4), 323–355.
- Tippett, L. (1931). *The Methods of Statistics* (1st ed.). London: Williams and Norgate.

- West, B., K. Welch, and A. Galecki (2007). *Linear Mixed Models: a Practical Guide Using Statistical Software*. Chapman and Hall CRC.
- Zewotir, T. and J. Galpin (2005). Influence diagnostics for linear mixed models. *Journal of Data Science* 3, 153–177.