

Contents

1	Introduction to Method Comparison Studies	2
1.1	Introduction	3
1.2	Grubbs' FCT Data	6
2	Bland-Altman Methodology	7
2.1	Bland-Altman Plots	8
2.2	Limits of Agreement	9
2.2.1	Prevalence of the Bland-Altman plot	10
2.3	Coefficient of Repeatability	11
2.4	Bartko's Ellipse	12
3	Regression Techniques	14
3.1	Deming Regression	14
3.2	Bradley Blackwood	14
3.3	Deming Regression	15
4	LME Techniques	16
4.1	Linear Mixed Effects Models	17
4.1.1	What are LME Models?	17
4.1.2	Laird-Ware Notation	17
4.2	The Research of Carstensen et al	18
4.2.1	Using LME Models for Method Comparison	18
4.2.2	Computing LoAs with LMEs	18
4.3	VC Matrix Types	19
4.3.1	Identity Matrix	19
4.3.2	Symmetry	19
4.3.3	Compound Symmetry	19
4.4	Roy's Hypothesis Tests	19
4.5	Omnibus Test	19

Chapter 1

Introduction to Method Comparison Studies

1.1 Introduction

The problem of assessing the agreement between two or more methods of measurement is ubiquitous in scientific research, and is commonly referred to as a ‘method comparison study’. Published examples of method comparison studies can be found in disciplines as diverse as pharmacology (?), anaesthesia (?), and cardiac imaging methods (?).

To illustrate the characteristics of a typical method comparison study consider the data in Table I (?). 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, identified here by the labels ‘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 1.1: Velocity measurement from the three chronographs (Grubbs 1973).

An important aspect of these data is that all three methods of measurement are assumed to have an attended measurement error, and the velocities reported in Table 1.1 can not be assumed to be ‘true values’ in any absolute sense.

A method of measurement should ideally be both accurate and precise. ? describes agreement as being a broader term that contains both of those qualities. An accurate measurement method will give results close to the unknown ‘true value’. The precision of a method is indicated by how tightly measurements obtained under identical conditions are distributed around their mean measurement value. A precise and accurate method will yield results consistently close to the true value. Of course a method may be accurate, but not precise, if the average of its measurements is close to the true value, but those measurements are highly dispersed. Conversely a method that is not accurate 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 smaller 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 one would conclude that there is lack of agreement between the two methods.

The absence of inter-method bias by itself is not sufficient to establish whether two measurement methods agree. The two methods must also have

equivalent levels of precision. Should one method yield results considerably more variable than those of the other, they can not be considered to be in agreement. With this in mind a methodology is required that allows an analyst to estimate the inter-method bias, and to compare the precision of both methods of measurement.

1.2 Grubbs' FCT Data

Chapter 2

Bland-Altman Methodology

2.1 Bland-Altman Plots

2.2 Limits of Agreement

2.2.1 Prevalence of the Bland-Altman plot

?, which further develops the Bland-Altman methodology, was found to be the sixth most cited paper of all time by the ?. ? describes the rate at which prevalence of the Bland-Altman plot has developed in scientific literature. ? reviewed the use of Bland-Altman plots by examining all articles in the journal ‘Clinical Chemistry’ between 1995 and 2001. This study concluded that use of the Bland-Altman plot increased over the years, from 8% in 1995 to 14% in 1996, and 31-36% in 2002.

The Bland-Altman Plot has since become expected, and often obligatory, approach for presenting method comparison studies in many scientific journals (?). Furthermore ? recommend its use in papers pertaining to method comparison studies for the journal of the British Hypertension Society.

2.3 Coefficient of Repeatability

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

2.4 Bartko's Ellipse

As a complement to the Bland-Altman plot, ? proposes the use of a bivariate confidence ellipse, constructed for a predetermined level. ? provides the relevant calculations for the ellipse. This ellipse is intended as a visual guidelines for the scatter plot, for detecting outliers and to assess the within- and between-subject variances.

The minor axis relates to the between subject variability, whereas the major axis relates to the error mean square, with the ellipse depicting the size of both relative to each other. Consequently Bartko's ellipse provides a visual aid to determining the relationship between variances. If $\text{var}(a)$ is greater than $\text{var}(d)$, the orientation of the ellipse is horizontal. Conversely if $\text{var}(a)$ is less than $\text{var}(d)$, the orientation of the ellipse is vertical.

The Bland-Altman plot for the Grubbs data, complemented by Bartko's ellipse, is depicted in Figure 1.7. The fourth observation is shown to be outside the bounds of the ellipse, indicating that it is a potential outlier.

The limitations of using bivariate approaches to outlier detection in the Bland-Altman plot can demonstrated using Bartko's ellipse. A covariate is added to the 'F vs C' comparison that has a difference value equal to the inter-method bias, and an average value that markedly deviates from the rest of the average values in the comparison, i.e. 786. Table 1.8 depicts a 95% confidence ellipse for this manipulated data set. By inspection of the confidence interval, a conclusion would be reached that this extra covariate is an outlier, in spite of the fact that this observation is wholly consistent with the conclusion of the Bland-Altman plot.

Importantly, outlier classification must be informed by the logic of the data's formulation. In the Bland-Altman plot, the horizontal displacement of any observation is supported by two independent measurements. Any observation should not be considered an outlier on the basis of a noticeable horizontal displacement from the main cluster, as in the case with the extra covariate. Conversely, the fourth observation, from the original data set, should be considered an outlier, as it has a noticeable vertical displacement from the rest of the observations.

In classifying whether a observation from a univariate data set is an outlier, many formal tests are available, such as the Grubbs test for outliers. In assessing whether a covariate in a Bland-Altman plot is an outlier, this test is useful when applied to the case-wise difference values treated as a univariate data set. The null hypothesis of the Grubbs test procedure is the absence of any outliers in the data set. Conversely, the alternative hypotheses is that there is at least one outlier present.

The test statistic for the Grubbs test (G) is the largest absolute deviation from the sample mean divided by the standard deviation of the differences,

$$G = \max_{i=1, \dots, n} \frac{|d_i - \bar{d}|}{S_d}.$$

For the 'F vs C' comparison it is the fourth observation gives rise to the test

statistic, $G = 3.64$. The critical value is calculated using Student's t distribution and the sample size,

$$U = \frac{n-1}{\sqrt{n}} \sqrt{\frac{t_{\alpha/(2n), n-2}^2}{n-2 + t_{\alpha/(2n), n-2}^2}}.$$

For this test $U = 0.75$. The conclusion of this test is that the fourth observation in the 'F vs C' comparison is an outlier, with p -value = 0.003, according with the previous result using Bartko's ellipse.

Chapter 3

Regression Techniques

3.1 Deming Regression

3.2 Bradley Blackwood

3.3 Deming Regression

Chapter 4

LME Techniques

4.1 Linear Mixed Effects Models

4.1.1 What are LME Models?

4.1.2 Laird-Ware Notation

4.2 The Research of Carstensen et al

4.2.1 Using LME Models for Method Comparison

4.2.2 Computing LoAs with LMEs

4.3 VC Matrix Types

4.3.1 Identity Matrix

4.3.2 Symmetry

4.3.3 Compound Symmetry

4.4 Roy's Hypothesis Tests

4.5 Omnibus Test