# Method Comparison Studies in Practise

# **Bendix Carstensen**

Steno Diabetes Center, Denmark
& Department of Biostatistics, University of Copenhagen
bxc@steno.dk www.biostat.ku.dk/~bxc

28–30 November 2007 Dept. of Medical Epidemiology and Biostatistics, Karolinska Institutet, Stockholm

# Introduction to computing

# Wednesday 28 November 2007, afternoon

### Bendix Carstensen

Method Comparison Studies in Practise 28–30 November 2007 Dept. of Medical Epidemiology and Biostatistics, Karolinska Institutet, Stockholm

# **Course structure**

The course is both theoretical and practical, i.e. the aim is to convey a basic understanding of the problems in method comparison studies, but also to convey practical skills in handling the statistical analysis.

- **R** for data manipulation ad graphics.
- WinBUGS for estimation in non-linear variance component models.





# Analyses/plots in this course

- Scatter plots.
- ▶ Bland-Altman plots (y x vs. (x + y)/2)
- Limits of agreement.
- Models with constant bias.
- Models with linear bias.
- Conversion formulae between methods (single replicates)
- Plots of converison equations.
- Graphical reporting of variance components.

Introduction to computing

### 8/64

# Requirements

- **R** for data manipulation and graphics:
- Tinn-R convenience editior with syntax highlighting for R.
- nlme-package for variance component models
   constant bias.
- WinBUGS for fitting models with linear bias (non-linear variance component models, over-parametrized).

All of it works from within **R**.

Introduction to computing

### 9/64

# Functions in the MethComp package

4 broad categories of functions in MethComp:

- Graphical just exploring data.
- Data manipulation reshaping and changing. Simulation.
- Analysis function fitting models to data.
- Reporting functions displaying the results from analyses.

# **Graphical functions**

- BA.plot Makes a Bland-Altman plot of two methods from a data frame with method comparison data, and computes limits of agreement. The plotting etc is really done by a call to
- BlandAltman Draws a Bland-Altman plot and computes limits of agreement.
- plot.meth Plots all methods against all other, both as a scatter plot and as a Bland-Altman plot.
- bothlines Adds regression lines of y on x and vice versa to a scatter plot.

11/ 64

# Data manipulating functions

Introduction to computing

- make.repl Generates a repl column in a data frame with columns meth, item and y.
- perm.repl Randomly permutes replicates within (method,item) and assigns new replicate numbers.
- to.wide Transforms a data frame in the long form to the wide form.
- to.long Reverses the result of to.wide.
- tab.repl Tabulates replicates by methods and items.
- sim.meth Simulates a dataset from a method comparison experiment for given parameters for bias, exchangeability and variances.

12/64

# **Analysis functions**

Introduction to computing

- Deming Performs Deming regression, i.e. regression with errors in both variables.
- BA.est Estimates in the variance components models underlying the concept of limits of agreement, and returns the bias and the variance components. Assumes constant bias between methods.
- MethComp Estimates via BUGS in the general model with non-constant bias (and in the future) possibly non-constant standard deviations of the variance components.
   Produces a MethComp object.

# **Reporting functions**

These functions all take a MethComp object as input.

- print.MethComp Prints a table of conversion equation between methods analyzed, with prediction standard deviations. Also gives summaries of the posteriors for the parameters that constitute the conversion algorithms.
- plot.MethComp Plots the conversion lines between methods with prediction limits.
- plot.VarComp Plots smoothed posterior densities for the variance component estimates.

Introduction to computing

### 14/64

# Does it work?

You should get something reasonable out of this:

```
library(MethComp)
data(ox)
plot.meth(ox)
plot.meth(perm.repl(ox))
BA.plot(ox)
BA.est(ox)
BA.est(perm.repl(ox))
MethComp(ox,code.only=TRUE)
m1 <- MethComp(ox)
print(m1)
plot(m1)
plot.VarComp(m1)</pre>
```

- if it works we are ready for tomorrow!

Introduction to computing

15/64

# Any practical examples?

# **Comparing two methods with one measurment on each**

# Thursday 29 November 2007, morning

# **Bendix Carstensen**

Method Comparison Studies in Practise 28–30 November 2007 Dept. of Medical Epidemiology and Biostatistics, Karolinska Institutet, Stockholm

# **Comparing measurement methods**

General questions:

- Are results systematically different?
- Can one method safely be replaced by another?
- What is the size of measurement errors?
- Different centres use different methods of measurement: How can we convert from one method to another?

17/64









# Model in "Limits of agreement"

Methods  $m = 1, \ldots, M$ , applied to  $i = 1, \ldots, I$  individuals:

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

- Two-way analysis of variance model, with unequal variances in columns.
- Different variances are not identifiable without replicate measurements for M = 2 because the variances cannot be separated.

23/64

# Limits of agreement:

Unequal variances induce correlation between  $D_i$  and  $A_i$ :

$$\operatorname{cov}(D_i, A_i) = \frac{1}{2}(\sigma_x^2 - \sigma_y^2) \neq 0 \quad \text{if } \sigma_x \neq \sigma_y$$

In correlation terms:

$$\rho(D,A) = \frac{1}{2} \frac{\sigma_x^2 - \sigma_y^2}{\sigma_x^2 + \sigma_y^2}$$

i.e. the correlation depends on whether the difference between the variances is large relative to the sizes of the two.

Models

Models

# Limits of agreement:

Usually interpreted as the likely differnence between two future measurements, one with each method:

$$\widehat{y_2 - y_1} = \hat{D} = \alpha_2 - \alpha_1 \pm 1.96 \, \text{s.d.}(D)$$

But it can of course also be converted to a prediction interval for  $y_2$  given  $y_1$ :

$$\hat{y}_2|y_1 = \alpha_2 - \alpha_1 + y_1 \pm 1.96 \,\mathrm{s.d.}(D)$$

### 25/64

# Repeatability and reproducibility

Models

# Thursday 29 November 2007, morning

# Bendix Carstensen

Method Comparison Studies in Practise 28–30 November 2007 Dept. of Medical Epidemiology and Biostatistics, Karolinska Institutet, Stockholm

# Accuracy of a measurement method

Repeatability:

The accuracy of the method under exactly similar circumstances; i.e. the same lab, the same technician, and the same day. (**Repeata**bility conditions)

Reproducibility:

The accuracy of the method under comparable circumstances, i.e. the same machinery, the same kit, but possibly different days or laboratories or technicians. (**Reproduci**bility conditions)

# **Quantification of accuracy**

- Upper limit of a 95% confidence interval for the difference between two measurments.
- Suppose the variance of the measurement is  $\sigma^2$ :

 $\operatorname{var}(y_{mi1} - y_{mi2}) = 2\sigma^2$ 

i.e the standard error is  $\sqrt{2}\sigma$ , and a confidnece interval for the difference:

 $0 \pm 1.96 \times \sqrt{2}\sigma = 0 \pm 2.772\sigma \approx 2.8\sigma$ 

 This is called the reproducibility coefficient or simply the reproducibility. (The number 2.8 is used as a convenient approximation).

Repeatability and reproducibility

27/64

# **Quantification of accuracy**

- Where do we get the  $\sigma$ ?
- Repeat measurements on the same item (or even better) several items.
- The conditions under which the repeat (replicate) measurements are taken determines whether we are estimating repeatability or reproducibility.
- In larger experiments we must consider the exchangeability of the replicates — i.e. which replicates are done under (exactly) similar conditions and which are not.

Repeatability and reproducibility

28/ 64

# **Comparing two methods with replicate measurements**

# Thursday 29 November 2007, morning

# **Bendix Carstensen**

# Extension of the model: replicate measurements

 $y_{mir} = \alpha_m + \mu_i + c_{mi} + e_{mir}$ s.d. $(c_{mi}) = \tau_m$  — "matrix"-effect s.d. $(e_{mir}) = \sigma_m$  — measurement error

- Replicates within (m, i) is needed to separate  $\tau$  and  $\sigma$ .
- ► Even with replicates, the *τ*s are only estimable if *M* > 2.
- Still assumes that the difference between methods is constant.
- Assumes exchangeability of replicates.

Comparing two methods with replicate measurements

# 29/64

# Extension of the model: replicate measurements

 $y_{mir} = \alpha_m + \mu_i + a_{ir} + c_{mi} + e_{mir}$ s.d. $(a_{ri}) = \omega$  — between replicates s.d. $(c_{mi}) = \tau_m$  — "matrix"-effect s.d. $(e_{mir}) = \sigma_m$  — measurement error

- Still assumes that the difference between methods is constant.
- Replicates are *linked* between methods: a<sub>ir</sub> is common across methods, i.e. the first replicate on a person is made under similar conditions for all methods (i.e. at a specific day or the like).

30/64

# **Replicate measurements**

Comparing two methods with replicate measurements

Two approaches to limits of agreement with replicate measurements:

- 1. Take means over replicates within each method by item stratum.
- 2. Replicates within item are taken as items.





# Wrong or almost right

In the model the correct limits of agreement would be:

$$\alpha_1 - \alpha_2 \pm 1.96\sqrt{\tau_1^2 + \tau_2^2 + \sigma_1^2 + \sigma_2^2}$$

If we are using means of replicates to form the differences we have:

$$\bar{d}_{i} = \bar{y}_{1i} - \bar{y}_{2i} = \alpha_{1} - \alpha_{2} + \frac{\sum_{r} a_{ir}}{R_{1i}} - \frac{\sum_{r} a_{ir}}{R_{2i}} + c_{1i} - c_{2i} + \frac{\sum_{r} e_{1ir}}{R_{1i}} - \frac{\sum_{r} e_{2ir}}{R_{2i}}$$

The terms with  $a_{ir}$  are only relevant for linked replicates in which case  $R_{1i} = R_{2i}$  and therefore the term vanishes. Thus:

 $\operatorname{var}(\bar{d}_i) = \tau_1^2 + \tau_2^2 + \sigma_1^2 / R_{1i} + \sigma_2^2 / R_{2i} < \tau_1^2 + \tau_2^2 + \sigma_1^2 + \sigma_2^2$ 

so the limits of agreement calculated based on the means are much too narrow as prediction limits for differences between future *single* measurements.

38/ 64

# (Linked) replicates as items

Comparing two methods with replicate measurement

If replicates are taken as items, then the calculated differences are:

 $d_{ir} = y_{1ir} - y_{2ir} = \alpha_1 - \alpha_2 + c_{1i} - c_{2i} + e_{1ir} - e_{2ir}$ 

which has variance  $\tau_1^2 + \tau_2^2 + \sigma_1^2 + \sigma_2^2$ , and so gives the correct limits of agreement. However, the differences are not independent:

$$\operatorname{cov}(d_{ir}, d_{is}) = \tau_1^2 + \tau_2^2$$

Negligible if the residual variances are very large compared to the interaction, variance likely to be only slightly downwards biased.

Comparing two methods with replicate measurements

39/64

# Exchangeable replicates as items?

If replicates are exchangeable it is not clear how to produce the differences using replicates as items.

If replicates are paired at random (se the function perm.repl), the variance will still be correct using the model without the  $i \times r$  interaction term  $(a_{ir})$ :

$$\operatorname{var}(y_{1ir} - y_{2is}) = \tau_1^2 + \sigma_1^2 + \tau_2^2 + \sigma_2^2$$

Differences will be positively correlated within item:

$$\operatorname{cov}(y_{1ir} - y_{2is}, y_{1it} - y_{2iu}) = \tau_1^2 + \tau_2^2$$

— slight underestimate of the true variance.





# A general model

# Thursday 29 November 2007, morning

# **Bendix Carstensen**

# **Extension of the model:**

 $\begin{array}{lll} y_{mir} &=& \alpha_m + \mu_i + a_{ir} + c_{mi} + d_{mr} + e_{mir} \\ & \mathrm{s.d.}(a_{ir}) = \omega & - \mathrm{between\ replicates} \\ & \mathrm{s.d.}(c_{mi}) = \tau_m & - \mathrm{``matrix''}\mathrm{-effect} \\ & \mathrm{s.d.}(d_{mr}) = \nu_m & - m \times r \\ & \mathrm{s.d.}(e_{mir}) = \sigma_m & - \mathrm{measurement\ error} \end{array}$ 

Method, Item, Replicate

- ▶ 1 3-way interaction
- ▶ 3 2-way interactions

What part of the interactions should be systematic (fixed) and what part should be random?

A general model

# $\left(m,r ight)$ - between replicates within method

This effect has  $M\times R$  levels, usually a rather small number.

This effect will therefore normally be modelled as a fixed effect, but not necessarily with  $M \times R$  parameters, presumably fewer.

If replicates are times of sampling or analysis, we may consider different time trends for each method, e.g.

$$d_{mr} = \gamma_m t_r$$

A random  $m \times r$ -effect would be hard to interpret.

A general model

### 44/64

43/64

# (i, r) - between replicates within individual

Observations with same (i, r) — but different method — will be correlated.

Use if all methods are applied to each item at

- different times
- at different locations
- at different conditions

This means there is a minimal structure to replicates — they are linked.

There might be further structure, e.g. a systematic effect of a time.

# (m,i) - between methods within individual

This is what is often called a "matrix" effect.

Matrix in the chemical sense: The surrounding matter ("matrix") in which the stuff of interest is dissolved.

Represents random effects of items reacting differently on each measurement method.

Logical to require that the variance of these methods was allowed to differ between methods.

A general model

### 46/64

47/64

# Variance component model!

Note we do not consider the method by replicate interaction any more.

The model is a (standard) variance component model, where two of the variance components depend on method.

A general model

# Fitting the variance component model Complicated and counter-intuitive in R: > library( nlme ) > lme( y ~ meth + item, random = list( item = pdIdent(~meth - 1), repl = ~1), weights = varIdent(form = ~1 | meth), data = ox)

A general model

```
Random effects:
    Formula: ~meth - 1 | item
    Structure: Multiple of an Identity
             methCO methpulse
   StdDev: 2.928042 2.928042
    Formula: ~1 | repl %in% item
          (Intercept) Residual
   StdDev: 3.415692 2.224868
   Variance function:
    Structure: Different standard deviations per stratum
    Formula: ~1 | meth
    Parameter estimates:
         CO
               pulse
   1.000000 1.795365
   Number of Observations: 354
   Number of Groups:
             item repl %in% item
               61
                              177
A general model
                                                       49/64
```

# Tease out variances for later use?

Even worse.

Therefore it has been packaged in a function that calls lme and then tease out the relevant parameters.

> BA.est(ox)

\$bias CO pulse 0.000000 -2.470446

\$sd.s
 MxI.CO MxI.pulse IxR resid.CO resid.pulse
 2.928042 2.928042 3.415692 2.224868 3.994451
Warning message:

In pt(q, df, lower.tail, log.p) : NaNs produced

50/64

# **Unequal bias**

A general model

# Thursday 29 November 2007, afternoon

# Bendix Carstensen

# Extension with non-constant bias

 $y_{mir} = \alpha_m + \beta_m \mu_i + random \text{ effects}$ 

There is now a *scaling* between the methods.

Methods do not measure on the same scale — the relative scaling is *estimated*, between method 1 and 2 the scale is  $\beta_2/\beta_1$ .

Consequence: Multiplication of all measurements on one method by a fixed number does not change results of analysis:

The corresponding  $\beta$  is multiplied by the same factor as is the variance components for this method.

Unequal bias

### 51/64

# Variance components

All two-way interactions:

$$y_{mir} = \alpha_m + \beta_m(\mu_i + a_{ir}) + c_{mi} + d_{mr} + e_{mir}$$

The random effects  $c_{mi}$ ,  $d_{mr}$  and  $e_{mir}$  have variances specific for each method.

But  $a_{ir}$  does not depend on m — must be scaled to each of the metods by the corresponding  $\beta$ .

Implies that  $\omega = \text{s.d.}(a_{ir})$  is irrelevant — the scale is arbitrary. The relevant quantities are  $\beta_m \omega$  — the between replicate variation within item *as measured on the mth scale*.

Unequal bias

# Variance components

Method, Item, Replicate.

$$y_{mir} = \alpha_m + \beta_m(\mu_i + a_{ir}) + c_{mi} + d_{mr} + e_{mir}$$
  
s.d. $(c_{mi}) = \tau_m$ 

Matrix-effect: Each item reacts differently to each method.

If only two methods compared:  $\tau_1$  and  $\tau_2$  cannot be separated:

$$y_{mir} = \alpha_m + \beta_m(\mu_i + a_{ir} + c_{mi}) + d_{mr} + e_{mir}$$
  
s.d. $(c_{mi}) = \tau$ 

Unequal bias

52/64

# Variance components

Method, Item, Replicate.

$$y_{mir} = \alpha_m + \beta_m(\mu_i + a_{ir}) + c_{mi} + d_{mr} + e_{mir}$$
  
s.d. $(d_{mr}) = \nu_m$ 

Number of methods and replicates are normally small.

More likely to be included as a fixed effect, for example as specific effects of analysis day for each method.

54/64

55/64

# Variance components

Method, Item, Replicate.

$$y_{mir} = \alpha_m + \beta_m(\mu_i + a_{ir}) + c_{mi} + d_{mr} + e_{mir}$$
  
s.d. $(a_{ir}) = \omega$ 

Common across methods — must be scaled relative to the methods.

Included if replicates are linked across methods, e.g. if there is a sequence in the replicates.

The relevant quantities to reports are  $\beta_m \omega$  — the s.d. on the scale of the *m*th method.

Unequal bias

Unequal bias

# **Conversion between methods**

# Friday 30 November 2007, morning

# **Bendix Carstensen**

# Predicting method 2 from method 1

The random effects have expectation 0, so:

$$E(y_{20r}|y_{10r}) = \hat{y}_{20r} = \alpha_2 + \frac{\beta_2}{\beta_1}(y_{k0r} - \alpha_1)$$

Conversion between methods

57/64

$$y_{20r} = \alpha_2 + \frac{\beta_2}{\beta_1} (y_{10r} - \alpha_1 - c_{10} - e_{10r}) + c_{20} + e_{20r}$$

$$\operatorname{var}(\hat{y}_{20r}|y_{10r}) = \left(\frac{\beta_2}{\beta_1}\right)^2 (\tau_1^2 + \sigma_1^2) + (\tau_2^2 + \sigma_2^2)$$

The slope of the prediction line from method 1 to method 2 is  $\beta_2/\beta_1.$ 

The width of the prediction interval is:

$$2 \times 1.96 \times \sqrt{\left(\frac{\beta_2}{\beta_1}\right)^2 (\tau_1^2 + \sigma_1^2) + (\tau_2^2 + \sigma_2^2)}$$

Conversion between methods

If we do the prediction the other way round  $(y_1|y_2)$  we get the same relationship i.e. a line with the inverse slope,  $\beta_1/\beta_2$ .

The width of the prediction interval in this direction is:

$$2 \times 1.96 \times \sqrt{(\tau_1^2 + \sigma_1^2) + \left(\frac{\beta_1}{\beta_2}\right)^2 (\tau_2^2 + \sigma_2^2)}$$
$$= 2 \times 1.96 \times \frac{\beta_1}{\beta_2} \sqrt{\left(\frac{\beta_2}{\beta_1}\right)^2 (\tau_1^2 + \sigma_1^2) + (\tau_2^2 + \sigma_2^2)}$$

i.e. if we draw the prediction limits as straight lines they can be used both ways.





Variance components

 $y_{mir} = \alpha_m + \beta_m(\mu_i + a_{ir}) + c_{mir} + e_{mir}$ 

The total variance of a measurement is:

$$\sqrt{\beta_m^2 \omega^2 + \tau_m^2 + \sigma_m^2}$$

These are the variance components reported by print.MethComp and shown by plot.VarComp.

62/64

# Repeatabiliy and reproducibility

Repeatability is based on the difference between measurements made under comparable, though not exactly identical conditions.

Reproducibility is based on the difference between measurements made under comparable, though not exactly identical conditions.

This is a different setting from the one underlying the modelling of data from a comparison experiment.

The exchangeability has no meaning, we are discussing future measurements in different circumstances.

Variance components

# **Repeatabiliy and reproducibility**

Repeatability:  $2.8\sigma_m$ :

same individual, same replicate, but not considering the variation that constitute differences between replicates *in the experiment*.

Hence *reproducibility* is not estimable from a classical experiment, unless an extra layer of replication is introduced — i.e. different laboratories.

Variance components

64/64