Analysis of Method Comparison Studies

Bendix Carstensen& Lyle Gurrin

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

MEGA center, School of Population Health University of Melbourne lgurrin@unimelb.edu.au www.epi.unimelb.edu.au/about/staff/gurrin-lyle

15 February 2008 MEGA Center, SPH, University of Melbourne Comparing two methods with one measurement on each

Friday 15 February

Bendix Carstensen

Analysis of Method Comparison Studies 15 February 2008 MEGA Center, SPH, University of Melbourne

(Comp-simple)

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?

Two methods for measuring fat content in human milk:



Comparing two methods with one measurement on each

Two methods — one measurement by each

How large is the difference between a measurement with method 1 and one with method 2 on a (randomly chosen) person?

$$D_i = y_{1i} - y_{2i}, \qquad D, \qquad \text{s.d.}(D)$$

"Limits of agreement:"

$$\bar{D} \pm 2 \times \text{s.d.}(D)$$

95% prediction interval for the difference between a measurement by method 1 and one by method 2. [?, ?]

Limits of agreement: Interpretation

- If a new patient is measured once with each of the two methods, the difference between the two values will with 95% probability be within the limits of agreement.
- This is a prediction interval for a (future) difference.
- Requires a clinical input: Are the limits of agreement sufficiently narrow to make the use of either of the methods clinically acceptable?
- Is it relevant to test if the mean is 0?

Limits of agreement: Test?

Testing whether the difference is 0 is a bad idea:

- If the study is sufficiently small this will be accepted even if the difference is important.
- If the study is sufficiently large this will be rejected even if the difference is clinically irrelevant.
- It is an equivalence problem: Clinical input is required!

Limits of agreement:



Plot differences (D_i) versus averages (A_i) .

Model in "Limits of agreement"

Methods m = 1, ..., M, applied to i = 1, ..., 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.

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.

Limits of agreement:

Usually interpreted as the likely difference 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 :

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

Introduction to computing

Friday 15 February

Bendix Carstensen

Analysis of Method Comparison Studies 15 February 2008 MEGA Center, SPH, University of Melbourne

(Intro-comp)

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 and graphics.
- WinBUGS for estimation in non-linear variance component models.

Software considerations

- R, SAS and Stata all have interfaces to WinBUGS.
- ▶ But **R** have more flexible graphical facilities.
- ► The MethComp package is written for **R**.

Therefore we use ${\boldsymbol{\mathsf{R}}}$ in this course.



Introduction to computing



Linked replicates.



Linked replicates.

Introduction to computing

How it works

Example data sets are included in the MethComp package. Contains the following variables.

meth — method
item — item, person, individual, sample
repl — replicate (if present)
y — the actual measurement

 — or rather *should* in order for the functions in MethComp to work.

How it looks

> sı	<pre>subset(ox,item<3)</pre>								
	meth	У							
1	CO	1	1	78.0					
2	CO	1	2	76.4					
3	CO	1	3	77.2					
4	CO	2	1	68.7					
5	CO	2	2	67.6					
6	CO	2	3	68.3					
184	pulse	1	1	71.0					
185	pulse	1	2	72.0					
186	pulse	1	3	73.0					
187	pulse	2	1	68.0					
188	pulse	2	2	67.0					
189	pulse	2	3	68.0					

<pre>> subset(to.wide(ox),item<3)</pre>									
Not	te:								
Replicate measurements are ta									
1	item	repl	id	CO	pulse				
1	1	1	1.1	78.0	71				
2	1	2	1.2	76.4	72				
3	1	3	1.3	77.2	73				
4	2	1	2.1	68.7	68				
5	2	2	2.2	67.6	67				
6	2	3	2.3	68.3	68				

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.

Requirements

- **R** for data manipulation and graphics:
- Tinn-R convenience editor 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 \mathbf{R} .

Functions in the MethComp package

5 broad categories of functions in MethComp:

- Graphical just exploring data.
- Data manipulation reshaping and changing.
- Simulation generating datasets.
- 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.

Data manipulating functions

- 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.

Analysis functions

- 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.

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)
```

Repeatability and reproducibility

Friday 15 February

Bendix Carstensen

Analysis of Method Comparison Studies 15 February 2008 MEGA Center, SPH, University of Melbourne

(Repro)

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 σ^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).

Quantification of accuracy

- Where do we get the σ ?
- 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.

Comparing two methods with replicate measurements

Friday 15 February

Bendix Carstensen

Analysis of Method Comparison Studies 15 February 2008 MEGA Center, SPH, University of Melbourne

(Comp-repl)

Extension of the model: replicate measurements

$$\begin{array}{lll} y_{mir} &=& \alpha_m + \mu_i + c_{mi} + e_{mir} \\ & & \mathrm{s.d.}(c_{mi}) = \tau_m & - \text{``matrix''-effect} \\ & & \mathrm{s.d.}(e_{mir}) = \sigma_m & - \text{measurement error} \end{array}$$

- Replicates within (m, i) is needed to separate τ and σ .
- ► 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.

Extension of the model: replicate measurements $y_{mir} = \alpha_m + \mu_i + a_{ir} + c_{mi} + e_{mir}$ s.d. $(a_{ir}) = \omega$ — 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_{ir} 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).

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.



Comparing two methods with replicate measurements






Replicate measurements

- The limits of agreement should still be for difference between future single measurements.
- Analysis based on the means of replicates is therefore wrong:
- Model:

$$y_{mir} = \alpha_m + \mu_i + a_{ir} + c_{mi} + e_{mir}$$
$$var(y_{1jr} - y_{2jr}) = \tau_1^2 + \tau_2^2 + \sigma_1^2 + \sigma_2^2$$
$$- note that the term a_{ir} - a_{ir} cancels because we are referring to the same replicate.$$

Recommendations

- Fit the correct model, and get the estimates from that, e.g. by using BA.est.
- If you must:
 - Use linked replicates as item.
 - If replicates are not linked; make a random linking. Note: If this give a substantially different picture than using the original replicate numbering as linking key, there might be something fishy about the data.



A general model

Friday 15 February

Bendix Carstensen

Analysis of Method Comparison Studies 15 February 2008 MEGA Center, SPH, University of Melbourne

(General)

Extension of the model:

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?

$\left(m,r ight)$ - 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.

$\left(i,r ight)$ - 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.

$\left(m,i ight)$ - 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.

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.

Fitting the variance component model

Complicated and counter-intuitive in R:

```
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
```

Tease out variances for later use?

Even worse.

Therefore it has been packaged in a function that calls 1me and then tease out the relevant parameters.

```
> BA.est(ox)
$bias
             pulse
      CO
0.00000 - 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
```

Unequal bias

Friday 15 February

Bendix Carstensen

Analysis of Method Comparison Studies 15 February 2008 MEGA Center, SPH, University of Melbourne

(Lin-bias)

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 β_2/β_1 .

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

The corresponding β is multiplied by the same factor as is the variance components for this method.

Two-way interactions:

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

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 methods by the corresponding β .

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 *m*th scale.

Method, Item, Replicate.

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

s.d. $(c_{mi}) = \tau_m$

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

If only two methods compared: τ_1 and τ_2 cannot be separated:

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

s.d.(c_{mi}) = τ

Method, Item, Replicate.

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

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.

Conversion between methods

Friday 15 February

Bendix Carstensen

Analysis of Method Comparison Studies 15 February 2008 MEGA Center, SPH, University of Melbourne

(Convert)

Predicting method 2 from method 1

The random effects have expectation 0, so:

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

$$y_{20r} = \alpha_2 + \frac{\beta_2}{\beta_1} (y_{10r} - \alpha_1 - c_{10} - e_{10r}) + c_{20} + e_{20r}$$
$$var(\hat{y}_{20r}|y_{10r}) = \left(\frac{\beta_2}{\beta_1}\right)^2 (\tau_1^2 + \sigma_1^2) + (\tau_2^2 + \sigma_2^2)$$

The slope of the prediction line from method 1 to method 2 is β_2/β_1 .

The width of the prediction interval is:

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

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, β_1/β_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.



Conversion between methods

What happened to the curvature?



Usually the prediction limits are curved:

$$\hat{y}|x \pm 1.96 \times \hat{\sigma}\sqrt{1 + x'x}$$

In our prediction we have ignored the last term (x'x), i.e. effectively assuming that there is no estimation error on $\alpha_{2\cdot 1}$ and $\beta_{2\cdot 1}$.

Friday 15 February

Bendix Carstensen

Analysis of Method Comparison Studies 15 February 2008 MEGA Center, SPH, University of Melbourne

(Var-comp)

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

3 variance components / random effects:

- *a_{ir}*: between replicates within item, ω²
 β_mω is the relevant quantity.
- c_{mi} : matrix effect τ_m^2 τ_m is the relevant quantity.
- e_{mir} : measurement error, residual variation σ_m^2 σ_m is the relevant quantity.

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

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.

Repeatability 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.

Repeatability 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.

Implementation in BUGS

Friday 15 February

Bendix Carstensen

Analysis of Method Comparison Studies 15 February 2008 MEGA Center, SPH, University of Melbourne

(BUGS-impl)

Implementation in BUGS

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

Non-linear hierarchical model: Implement in BUGS.

- ▶ The model is *symmetrical* in methods.
- Mean is overparametrized.
- Choose a prior (and hence posterior!) for the µs with finite support.
- Keeps the chains nicely in place.

Results from fitting the model

The posterior dist'n of $(\alpha_m, \beta_m, \mu_i)$ is singular.

But the relevant translation quantities are identifiable:

$$\alpha_{2 \cdot 1} = \alpha_2 - \alpha_1 \beta_2 / \beta_1$$
$$\beta_{2 \cdot 1} = \beta_2 / \beta_1$$

So are the variance components.

Posterior medians used to devise prediction equations with limits.














Morale

- ► Use a proper model for your problem.
- Get the exchangeability right.
- ▶ Report the model in a useful way.

The MethComp package for R

Implemented model:

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

- Replicates required.
- R2WinBUGS is required.
- Dataframe with variables meth, item, repl and y.
- The function MethComp writes a BUGS-program, initial values and data to files.
- Runs WinBUGS and sucks results back in to R, and gives a nice overview of the conversion equations.

Example output: Oximetry

> ox.mi.ir <- MethComp(ox, n.iter=5000)
> ox.mi.ri

Comparison of 2 methods, using 354 measurements on 61 individuals, with up to 3 replicate measurements. (2 * 61 * 3 = 366):

No. individuals with measurements on each method: # replicates Method 1 2 3 Sum CO 1 4 56 61 pulse 1 4 56 61

Example output: Oximetry

Conv	ersion	formu	ılae (y ₋	to = alph	ia + beta*y	y_from	+/- 2*sd	.pred):
	From:	CO			pulse			
		alpha	beta	sd.pred	alpha	beta	sd.pred	
To:								
CO		0.000	1.000	4.266	-4.328	1.098	8.487	
puls	е	3.939	0.911	7.606	0.000	1.000	5.534	

Variance components (standard deviations):

	50%	2.5%	97.5%	0%	100%
sigma.mir[CO]	1.6285	0.2092	2.8274	0.0724	3.4330
<pre>sigma.mir[pulse]</pre>	4.2580	3.5390	4.9725	3.0670	5.9800
sigma.mi[CO]	4.8043	2.7504	13.3685	2.2597	17.6134
sigma.mi[pulse]	4.3123	2.4981	11.5859	1.9248	13.2186
sigma.ir[CO]	3.9213	3.1452	4.7038	2.7289	5.3129
sigma.ir[pulse]	3.5433	2.7542	4.3516	2.2610	4.8723

HbA_{1c} - 3 different instruments

```
> hbv.mi.ir <- MethComp( hbv, n.iter=5000 )
> print( hbv.mi.ir, across=FALSE )
```

```
Conversion formula:
y_to = alpha + beta * y_from +/- 2*sd.pred:
          From: BR.V2 BR.VC Tosoh
To:
BR.V2 alpha
                 0.000 - 1.627 1.413
     beta
                1.000 1.154 0.946
     sd.pred
                0.254 2.079 2.099
BR.VC alpha
                1.417 0.000 2.412
     beta
                 0.867 1.000 0.819
     sd.pred 1.800 0.164 1.927
Tosoh alpha -1.591 -3.144 0.000
            1.057 1.220 1.000
     beta
                 2.145 2.249 0.156
     sd.pred
```

HbA_{1c} - 3 different instruments

Variance components (standard deviations):

50%2.5%97.5%0%100%sigma.mir[BR.V2]0.20890.18160.24010.16140.2692sigma.mir[BR.VC]0.10740.08130.12860.06420.1467sigma.mir[Tosoh]0.03450.00060.08240.00040.0984sigma.mi[BR.V2]1.34951.07801.77420.91942.1615sigma.mi[BR.VC]1.30881.04981.69790.86152.1350sigma.mi[Tosoh]1.44161.07825.36530.92506.3534sigma.ir[BR.V2]0.14180.10370.18820.08550.2319sigma.ir[BR.VC]0.12390.09280.15720.07970.1827sigma.ir[Tosoh]0.14960.12310.18150.09500.2002











The MethComp package

Also (currently) contains:

- BA.plot make a Bland-Altman plot and compute limits of agreement.
- BA.est estimates in the variance component model for the constant bias situation.
- Deming regression with errors in both variables.

A .pdf with a detailed derivation of the formulae (by Anders C Jensen) is included in the package too.

 A number of example data sets, amongst them all examples from [?].