# Analysis of Longitudinal Data in Stata, Splus and SAS

Rino Bellocco, Sc.D.
Department of Medical Epidemiology
Karolinska Institutet
Stockholm, Sweden

`rino@mep.ki.se`

March 12, 2001

# OUTLINE

- Longitudinal data

- Review

- Sample data set

- STATA (XTGEE, XTREG, GLLAMM6)

- SAS (Proc Mixed (Repetead, Random), Proc Glinmix,Proc Genmod)

- Splus (LME, YAGS)

- References

# Longitudinal Data

- **Longitudinal Studies:** studies in which the outcome variable is measured repeatedly over time. We do not necessarily require the same number of observations on each subject or that measurements be taken at the same times.

  $y_{ij} =$ value of j$^{th}$ observation on the $i^{th}$ subject

  measures at time $t_{ij}$.

- **Repeated measures:** Older term used for a special set of longitudinal designs with measurements at a common set of occasions, usually in an experimental design.

- Models for the analysis of longitudinal data can be considered a special case of generalized linear models, with the peculiar feature that the residuals terms are correlated, as the observations at different time points in a longitudinal study are taken on the same subject. Any of the model being proposed must take this dependence into account.

# Potential Advantages of Longitudinal Studies

- Allow investigation of events that occur in time; essential to the study of normal growth and ageing.

- Essential to the study of temporal patterns of response to treatments.

- Permit more complete ascertainment of exposure histories in epidemiological studies.

- Reduce unexplained variability in the response by using subject as his or her own control.

# Normally Distributed Data - Marginal Models

With longitudinal data, we can consider models of the form

$$Y_{ij} = \beta_0 + \beta_1 X_{1ij} + \beta_2 X_{2ij} + \ldots + \beta_Q X_{Qij} + \epsilon_{ij}$$

where the $\epsilon_{ij}$ are correlated within individuals (i.e. $Cov(\epsilon_{ij}, \epsilon_{ik}) \neq 0$) and the covariates $(X_{1ij}, ..., X_{Qij})$ include time, $t_{ij}$ (or indicators of time trends), treatment/exposure indicators and their interactions.

Recall that the "compound symmetry" assumption is unrealistic for longitudinal studies, instead we need to consider alternative models for $Cov(\epsilon_{ij}, \epsilon_{ik})$.

# Models for the Covariance

:

Note that with $p$ repeated measures, there are $\frac{p(p+1)}{2}$ parameters in the covariance matrix.

In selecting a model for the covariance matrix, a balance must be struck:

- With too little structure (e.g., unstructured). there may be too many parameters to be estimated with a limited amount of data (information) available $\implies$ weaker inferences concerning $\beta$

- With too much structure (e.g., compound symmetry), there is more information available for estimating $\beta$ but the potential risk of model misspecification $\implies$ apparently stronger, but potentially biased, inferences concerning $\beta$

# Other models

A number of additional models for the covariance that may be suitable for longitudinal data are

1. Autoregressive: The first-order autoregressive model, AR(1), has covariances of the form,
   $Cov(Y_{ij}, Y_{ik}) = \sigma^2 \rho^{|j-k|}$,
   i.e., homogeneous variances and correlations that decline over time.

$$
\begin{array}{cc}
 & \text{occasion} \\
\text{occasion} &
\begin{array}{c}
\phantom{1} \\ 1 \\ 2 \\ 3 \\ 4
\end{array}
\begin{array}{cccc}
1 & 2 & 3 & 4 \\
\left[\begin{array}{cccc}
1 & \rho & \rho^2 & \rho^3 \\
\rho & 1 & \rho & \rho^2 \\
\rho^2 & \rho & 1 & \rho \\
\rho^3 & \rho^2 & \rho & 1
\end{array}\right]
\end{array}
\end{array}
$$

   Autoregressive models are appropriate for equally-spaced measurement.

2. Exponential correlation models can handle unequally-spaced measurements.

   Suppose that measurements are made at times $t_j$, then the covariances are of the form,

$$
Cov(Y_{ij}, Y_{ik}) = \sigma^2 \rho^{|t_j - t_k|}.
$$

# STATA

**xtgee** fits generalized linear models of $Y_{ij}$, with covariates $X_{ij}$. Main components of a model:

1. **family** - assumed distribution of the response variables

2. **link** - link between response and its linear predictor

3. **corr** - structure of the working correlation

# Stata-xtgee

```
************************************************
*
* Sample program for NASUG 2001
* Data set: depress.dat from Hasbekt & Everitt
* Rino Bellocco
************************************************
infile subj group pre dep1 dep2 dep3 dep4 dep5 dep6
using c:\rino\nasug\depress.dat, clear
(61 observations read)


subj  group  pre dep1  dep2  dep3 dep4  dep5 dep6
 1     0     18  17    18    15   17    14   15
 2     0     27  26    23    18   17    12   10


Observations are correlated!

     | pre  dep1  dep2  dep3  dep4   dep5   dep6
-----+-------------------------------------------------------
 pre | 1.0000
dep1 | 0.2027   1.0000
dep2 | 0.2292   0.1937   1.0000
dep3 | 0.1683   0.0700   0.5645   1.0000
dep4 | 0.0561   0.0594   0.5125   0.9015   1.0000
dep5 | 0.1160   0.0654   0.5256   0.9160   0.9606   1.0000
dep6 | 0.1037   0.0184   0.5045   0.9035   0.9499   0.9743  1.0000
```

# Stata-xtgee

First step is to reshape the data so that we can use models.

reshape long dep, i(subj) j(visit) (note: j = 1 2 3 4 5 6)

| subj | visit | group | pre | dep |
|------|-------|-------|-----|-----|
| 1 | 1 | 0 | 18 | 17 |
| 1 | 2 | 0 | 18 | 18 |
| 1 | 3 | 0 | 18 | 15 |
| 1 | 4 | 0 | 18 | 17 |
| 1 | 5 | 0 | 18 | 14 |
| 1 | 6 | 0 | 18 | 15 |
| 2 | 1 | 0 | 27 | 26 |
| 2 | 2 | 0 | 27 | 23 |
| 2 | 3 | 0 | 27 | 18 |
| 2 | 4 | 0 | 27 | 17 |
| 2 | 5 | 0 | 27 | 12 |
| 2 | 6 | 0 | 27 | 10 |

# Stata-xtgee

## First, I run a model with independence structure

```
xtgee dep group pre visit, i(subj) t(visit) corr(indep) link(iden) fam(normal) nmp
```

```
GEE population-averaged model          Number of obs      =        295
Group variable:                  subj  Number of groups   =         61
Link:                        identity  Obs per group: min =          1
Family:                      Gaussian                 avg =        4.8
Correlation:              independent                 max =          6
                                       Wald chi2(3)       =     144.15
Scale parameter:              25.80052  Prob > chi2       =     0.0000

Pearson chi2(291):             7507.95  Deviance          =    7507.95
Dispersion (Pearson):         25.80052  Dispersion        =   25.80052
```

```
------------------------------------------------------------------------------
         dep |      Coef.   Std. Err.      z    P>|z|     [95% Conf. Interval]
-------------+----------------------------------------------------------------
       group |  -4.290664   .6072954    -7.07   0.000    -5.480941   -3.100387
         pre |   .4769071   .0798565     5.97   0.000     .3203913     .633423
       visit |  -1.307841    .169842    -7.70   0.000    -1.640725   -.9749569
       _cons |   8.233577   1.803945     4.56   0.000     4.697909    11.76924
------------------------------------------------------------------------------
```

# Stata-xtgee

## Then I fit a GLM with an exchangeable structure

```
. xtgee dep group pre visit, i(subj) t(visit) corr(exc) link(iden) fam(normal)

Iteration 1: tolerance = .04984936
Iteration 2: tolerance = .0004433
Iteration 3: tolerance = 4.602e-06
Iteration 4: tolerance = 4.782e-08

GEE population-averaged model          Number of obs      =         295
Group variable:                 subj   Number of groups   =          61
Link:                       identity   Obs per group: min =           1
Family:                     Gaussian                  avg =         4.8
Correlation:            exchangeable                  max =           6
                                       Wald chi2(3)       =      135.08
Scale parameter:             25.56569  Prob > chi2        =      0.0000

------------------------------------------------------------------------------
         dep |      Coef.   Std. Err.      z    P>|z|     [95% Conf. Interval]
-------------+----------------------------------------------------------------
       group |  -4.024676   1.081131    -3.72   0.000    -6.143654   -1.905698
         pre |   .4599018   .1441533     3.19   0.001     .1773666    .742437
       visit |  -1.226764   .1175009   -10.44   0.000    -1.457062   -.9964666
       _cons |   8.432806   3.120987     2.70   0.007     2.315783    14.54983
------------------------------------------------------------------------------
```

# Stata-xtgee

## Then I fit a model with unstructured correlation

```
xtgee dep group pre visit, i(subj) t(visit) corr(uns) link(iden) fam(normal)

GEE population-averaged model          Number of obs      =        295
Group and time vars:        subj visit  Number of groups   =         61
Link:                         identity  Obs per group: min =          1
Family:                       Gaussian                 avg =        4.8
Correlation:              unstructured                 max =          6
                                        Wald chi2(3)       =      94.13
Scale parameter:               25.87029  Prob > chi2       =     0.0000


------------------------------------------------------------------------------
        dep |     Coef.   Std. Err.      z    P>|z|    [95% Conf. Interval]
------------+-----------------------------------------------------------------
      group | -4.134413   .9986306    -4.14   0.000   -6.091693   -2.177133
        pre |  .3399185   .1326684     2.56   0.010    .0798932    .5999437
      visit | -1.228327   .1492831    -8.23   0.000   -1.520916   -.9357372
      _cons |  11.13045   2.892903     3.85   0.000    5.460464    16.80044
------------------------------------------------------------------------------
```

# Stata-xtgee

## And finally a model with AR1 structure

```
xtgee dep group pre visit, i(subj) t(visit) corr(ar1) link(iden) fam(normal)
note:  some groups have fewer than 2 observations
       not possible to estimate correlations for those groups
       8 groups omitted from estimation



Iteration 1: tolerance = .10070858
Iteration 2: tolerance = .00136623
Iteration 3: tolerance = .00002736
Iteration 4: tolerance = 5.508e-07

GEE population-averaged model          Number of obs     =        287
Group and time vars:        subj visit Number of groups  =         53
Link:                         identity Obs per group: min =          2
Family:                       Gaussian                avg =        5.4
Correlation:                     AR(1)                max =          6
                                       Wald chi2(3)      =      64.55
Scale parameter:              25.82413 Prob > chi2       =     0.0000


------------------------------------------------------------------------------
        dep |     Coef.   Std. Err.      z    P>|z|     [95% Conf. Interval]
-------------+----------------------------------------------------------------
      group |  -4.218194   1.053504    -4.00   0.000    -6.283023   -2.153364
        pre |   .4268002   .1376156     3.10   0.002     .1570785    .6965219
      visit |  -1.181975   .1907298    -6.20   0.000    -1.555799   -.8081517
      _cons |   9.037864   3.036076     2.98   0.003     3.087264    14.98846
------------------------------------------------------------------------------
```

# SAS-GLM

Here, I show what I think is the equivalent procedure in SAS (codes are reported at the end). Independence:

```
The REG Procedure
Model: MODEL1
Dependent Variable: dep

Analysis of Variance
```

|  |  | Sum of | Mean |  |  |
| Source | DF | Squares | Square | F Value | Pr > F |
| Model | 3 | 3719.12937 | 1239.70979 | 48.05 | <.0001 |
| Error | 291 | 7507.95172 | 25.80052 |  |  |
| Corrected Total | 294 | 11227 |  |  |  |

| Root MSE | 5.07942 | R-Square | 0.3313 |
| Dependent Mean | 11.32915 | Adj R-Sq | 0.3244 |
| Coeff Var | 44.83496 |  |  |

```
Parameter Estimates
```

|  |  | Parameter | Standard |  |  |
| Variable | DF | Estimate | Error | t Value | Pr > |t| |
| Intercept | 1 | 8.23358 | 1.80395 | 4.56 | <.0001 |
| group | 1 | -4.29066 | 0.60730 | -7.07 | <.0001 |
| pre | 1 | 0.47691 | 0.07986 | 5.97 | <.0001 |
| visit | 1 | -1.30784 | 0.16984 | -7.70 | <.0001 |

# SAS-GLM

Unrestricted Covariance structure

| Effect | group | Estimate | Standard Error | DF | t Value | Pr > \|t\| |
|---|---|---|---|---|---|---|
| Intercept | | 6.2422 | 2.8737 | 58 | 2.17 | 0.0339 |
| group | 0 | 4.1207 | 0.9739 | 58 | 4.23 | <.0001 |
| group | 1 | 0 | . | . | . | . |
| pre | | 0.3641 | 0.1292 | 58 | 2.82 | 0.0066 |
| visit | | -1.1091 | 0.1426 | 58 | -7.78 | <.0001 |

Compound structure

| Effect | group | Estimate | Standard Error | DF | t Value | Pr > \|t\| |
|---|---|---|---|---|---|---|
| Intercept | | 4.4124 | 3.1901 | 58 | 1.38 | 0.1719 |
| group | 0 | 4.0216 | 1.0887 | 58 | 3.69 | 0.0005 |
| group | 1 | 0 | . | . | . | . |
| pre | | 0.4598 | 0.1452 | 58 | 3.17 | 0.0025 |
| visit | | -1.2259 | 0.1167 | 233 | -10.50 | <.0001 |

AR1 structure

| Effect | group | Estimate | Error | DF | t Value | Pr > \|t\| |
|---|---|---|---|---|---|---|
| Intercept | | 5.0946 | 2.9691 | 58 | 1.72 | 0.0915 |
| group | 0 | 4.0317 | 1.0015 | 58 | 4.03 | 0.0002 |
| group | 1 | 0 | . | . | . | . |
| pre | | 0.4296 | 0.1331 | 58 | 3.23 | 0.0021 |
| visit | | -1.2221 | 0.1844 | 233 | -6.63 | <.0001 |

# SAS-GLM

```
libname rino 'c:\rino\nasug';
data rino;
infile  'c:\rino\nasug\depress.dat';
input subj group pre dep1 dep2 dep3 dep4 dep5 dep6;

if dep1=-9 then dep1=.     ;
if dep2=-9 then dep2=.     ;
if dep3=-9 then dep3=.     ;
if dep4=-9 then dep4=.     ;
if dep5=-9 then dep5=.     ;
if dep6=-9 then dep6=.     ;
run;
proc means;
var dep1 dep2 dep3 dep4 dep5 dep6 group pre;
run;

data rino1;
set rino;

visit=1; dep=dep1;t=1;output;
visit=2; dep=dep2;t=2;output;
visit=3; dep=dep3;t=3;output;
visit=4; dep=dep4;t=4;output;
visit=5; dep=dep5;t=5;output;
visit=6; dep=dep6;t=6;output;
run;
proc means;
var dep time pre group;
run;

/* proc print data=rino1;
   run;
*/
proc reg data=rino1;
model dep=group pre visit ;
run;

proc mixed data=rino1 noclprint method=ml ;
class subj group  t;
```

```
model dep = group pre visit /s;
repeated t /type=un subject=subj r;
title 'unrest.cov. structure, linear trend, ML';
run;

proc mixed data=rino1 noclprint method=ml;
class subj group t;
model dep = group  pre visit /s;
repeated t /type=cs subject=subj r;
title 'compound structure, linear trend, ML';
run;

proc mixed data=rino1 noclprint method=ml;
class subj group t;
model dep = group  pre visit /s;
repeated t /type=ar(1) subject=subj r;
title 'ar1 structure, linear trend, ML';
run;

proc mixed data=rino1 noclprint method=ml;
class subj group  t;
model dep = group  pre visit /s;
random intercept /type =un sub=subj s;
title 'random intercept, linear trend, ML';
run;

proc mixed data=rino1 noclprint method=ml;
class subj group  t;
model dep = group  pre visit /s;
random intercept visit /type =un sub=subj s;
title 'random intercept, linear trend, ML';
run;
```

# Stata SAS- comparison

Similar results are observed, however not the same estimates are produced. Testing and comparison of models with different covariance structures will be reported in a future paper (most likely an STB bullettin).

# Normally Distributed Data
# Random Effect Models

This approach assumes that the correlation arises among repeated measures as the regression coefficients vary across individuals.

That is, each subject is assumed to have an (unobserved) underlying level of response which persists across the $p$ measurements.

This subject effect is treated as random and the model becomes

$$Y_{ij} = \beta_0 + \beta_1 X_{1ij} + \beta_2 X_{2ij} + \ldots + \beta_{p-1} X_{p-1,ij} + b_i + e_{ij}$$

or

$$Y_{ij} = (\beta_0 + b_i) + \beta_1 X_{1ij} + \beta_2 X_{2ij} + \ldots + \beta_{p-1} X_{p-1,ij} + e_{ij}$$

(also known as "random intercepts model").

In the model

$$Y_{ij} = \beta_0 + \beta_1 X_{1ij} + \beta_2 X_{2ij} + \ldots + \beta_{p-1} X_{p-1,ij} + b_i + e_{ij}$$

the response for the $i^{th}$ subject is assumed to differ from the population mean, by a subject effect, $b_i$, and a within-subject measurement error, $e_{ij}$.

Alternatively, we have decomposed

$$\epsilon_{ij} = b_i + e_{ij}.$$

Furthermore, it is assumed that

$$b_i \stackrel{d}{=} N(0, \sigma_b^2); \qquad e_{ij} \stackrel{d}{=} N(0, \sigma_e^2)$$

and that $b_i$ and $e_{ij}$ are mutually independent.

The introduction of a random subject effect induces correlation among the repeated measures.

It can be shown that the following correlation structure results:

$$\text{Var}(Y_{ij}) = \sigma_b^2 + \sigma_e^2$$

$$\text{Cov}(Y_{ij}, Y_{ik}) = \sigma_b^2$$

$$\implies \text{Corr}(Y_{ij}, Y_{lj}) = \frac{\sigma_b^2}{\sigma_b^2 + \sigma_e^2}$$

$=$ correlation of observations on the same individual

**Stata** can fit this model using the **XTREG** procedure.

# XTREG/Stata

```
. xtreg dep group pre visit, i(subj) mle
```

```
Random-effects ML regression              Number of obs      =        295
Group variable (i) : subj                 Number of groups   =         61

Random effects u_i ~ Gaussian             Obs per group: min =          1
                                                         avg =        4.8
                                                         max =          6

                                          LR chi2(3)         =     111.62
Log likelihood  = -832.36607              Prob > chi2        =     0.0000
------------------------------------------------------------------------------
         dep |      Coef.   Std. Err.      z    P>|z|     [95% Conf. Interval]
-------------+----------------------------------------------------------------
       group |  -4.021599    1.08894    -3.69   0.000    -6.155882   -1.887316
         pre |   .4597672   .1451952     3.17   0.002     .1751898    .7443446
       visit |  -1.225857   .1168668   -10.49   0.000    -1.454912   -.9968024
       _cons |   8.434001   3.142894     2.68   0.007     2.274042    14.59396
-------------+----------------------------------------------------------------
     /sigma_u |   3.805795   .4160801     9.15   0.000     2.990293    4.621297
     /sigma_e |   3.346938     .15434    21.69   0.000     3.044438    3.649439
-------------+----------------------------------------------------------------
         rho |   .5638883   .0600327                       .4451442    .6771015
------------------------------------------------------------------------------
Likelihood ratio test of sigma_u=0: chibar2(01)=  127.28
Prob>=chibar2 = 0.000
```

# SAS

random intercept, linear trend, ML

Covariance Parameter Estimates

| Cov Parm | Subject | Estimate |
|----------|---------|----------|
| UN(1,1) | subj | 14.4836 |
| Residual | | 11.2021 |

Fit Statistics

| | |
|---|---|
| -2 Log Likelihood | 1664.7 |
| AIC (smaller is better) | 1676.7 |
| AICC (smaller is better) | 1677.0 |
| BIC (smaller is better) | 1689.4 |

Solution for Fixed Effects

| Effect | group | Estimate | Standard Error | DF | t Value | Pr > \|t\| |
|--------|-------|----------|----------------|-----|---------|-----------|
| Intercept | | 4.4124 | 3.1901 | 58 | 1.38 | 0.1719 |
| group | 0 | 4.0216 | 1.0887 | 233 | 3.69 | 0.0003 |
| group | 1 | 0 | . | . | . | . |
| pre | | 0.4598 | 0.1452 | 233 | 3.17 | 0.0017 |
| visit | | -1.2259 | 0.1167 | 233 | -10.50 | <.0001 |

Type 3 Tests of Fixed Effects

| Effect | Num DF | Den DF | F Value | Pr > F |
|--------|--------|--------|---------|--------|
| group | 1 | 233 | 13.64 | 0.0003 |
| pre | 1 | 233 | 10.03 | 0.0017 |
| visit | 1 | 233 | 110.35 | <.0001 |

# Splus

```
> summary(rem0)
Linear mixed-effects model fit by REML
 Data: rino
       AIC       BIC     logLik
  1678.536 1700.576 -833.2679

Random effects:
 Formula: visit ~ 1 | subj
        (Intercept) Residual
StdDev:    3.923239 3.353891

Fixed effects: dep ~ visit + pre + group
               Value Std.Error  DF   t-value p-value
(Intercept) 8.435886  3.224813 233   2.61593  0.0095
      visit -1.224393  0.117018 233 -10.46327  <.0001
        pre  0.459552  0.149022  58   3.08379  0.0031
      group -4.016623  1.117115  58  -3.59553  0.0007

 Correlation:
      (Intr)  visit    pre
visit -0.107
  pre -0.960  0.005
group -0.130 -0.040 -0.066

Standardized Within-Group Residuals:
       Min         Q1         Med        Q3       Max
 -3.840718 -0.5559042 -0.03438542 0.4645086 3.912141

Number of Observations: 295 Number of Groups: 61
```

# Random Intercepts and Slopes Models

A natural extension of the random intercepts model. The introduction of random intercepts and slopes induces a covariance matrix that depends on time $(t_{ij})$.

Consider the following model with intercepts and slopes that vary randomly among subjects

$$Y_{ij} = \beta_0 + \beta_1 t_{ij} + b_{i0} + b_{i1} t_{ij} + e_{ij}$$

Assume that $b_{i0}$ and $b_{i1}$ have mean zero and let $Var(e_{ij}) = \sigma_e^2$, $Var(b_{i0}) = \sigma_{00}^2$, $Var(b_{i1}) = \sigma_{11}^2$, and $Cov(b_{i0}, b_{i1}) = \sigma_{01}$.

Then, it can be shown that

$$Var(Y_{ij}) = \sigma_{00}^2 + 2t_{ij}\sigma_{01} + \sigma_{11}^2 t_{ij}^2 + \sigma_e^2$$

and

$$Cov(Y_{ij}, Y_{ik}) = \sigma_{00}^2 + (t_{ij} + t_{ik})\sigma_{01} + \sigma_{11}^2 t_{ij} t_{ik}$$

That is, the covariance matrix is a function of time. Stata has limited resources for modeling longitudinal data (GLLAMM6 is a routine provided

by Rabe-Hesketh which allows to fits this model, but it is not part of regular Stata and as, Sophia has told me, GLLAMM6 is intended for non-normal data where no exact method exists; instead we can use PROC MIXED in SAS and LME in Splus.

# STATA

```
gen cons=1 eq cons: cons eq slope: visit
gllamm6 dep group pre
visit, i(subj) nrf(2) eqs(cons slope) trace

gllamm model

log likelihood = -820.90341
```

```
-----------------------------------------------------------------------------
        dep |     Coef.   Std. Err.      z    P>|z|     [95% Conf. Interval]
------------+----------------------------------------------------------------
      group |  -3.459758   .9574966    -3.61   0.000    -5.336417   -1.583099
        pre |   .5769432   .0954126     6.05   0.000     .3899379    .7639484
      visit |  -1.240965   .1552877    -7.99   0.000    -1.545324   -.9366072
      _cons |   5.499468   2.249447     2.44   0.014     1.090632    9.908304
-----------------------------------------------------------------------------
```

Variance at level 1
-----------------------------------------------------------------------------
  8.1725165 (.86878708)

Variances and covariances of random effects
-----------------------------------------------------------------------------

***level 2 (subj)

    var(1): 23.758474 (5.8717413)
    cov(1,2): -2.2504823 (.98450321) cor(1,2): -.53217727

    var(2): .75269674 (.18593369)
-----------------------------------------------------------------------------

# SAS

random intercept + slope, linear trend, ML

Iteration History

| Iteration | Evaluations | -2 Log Like | Criterion |
|---|---|---|---|
| 0 | 1 | 1792.01280464 | |
| 1 | 2 | 1642.82321420 | 0.00000252 |
| 2 | 1 | 1642.82181110 | 0.00000000 |

Convergence criteria met.

Covariance Parameter Estimates

| Cov Parm | Subject | Estimate |
|---|---|---|
| UN(1,1) | subj | 22.3135 |
| UN(2,1) | subj | -2.4981 |
| UN(2,2) | subj | 0.8352 |
| Residual | | 8.3660 |

Fit Statistics

| -2 Log Likelihood | 1642.8 |
|---|---|
| AIC (smaller is better) | 1658.8 |

random intercept + slope, linear trend, ML

12:30 Saturday, Mar

The Mixed Procedure

Fit Statistics

```
            AICC (smaller is better)        1659.3
            BIC (smaller is better)         1675.7


              Null Model Likelihood Ratio Test


              DF      Chi-Square      Pr > ChiSq


               3        149.19          <.0001



                 Solution for Fixed Effects


                              Standard
    Effect      group    Estimate      Error      DF    t Value    Pr > |t|

    Intercept              4.2101     3.2138      58       1.31      0.1954
    group        0         4.0397     1.0922     181       3.70      0.0003
    group        1              0         .        .          .          .
    pre                    0.4682     0.1456     181       3.22      0.0015
    visit                 -1.2097     0.1651      52      -7.33      <.0001
```

# Splus

```
> summary(rem1)
Linear mixed-effects model fit by REML
 Data: rino
       AIC      BIC    logLik
  1659.905 1689.292 -821.9527

Random effects:
 Formula:  ~ visit | subj
 Structure: General positive-definite
              StdDev    Corr
(Intercept) 4.8414891 (Inter
      visit 0.9303804 -0.572
   Residual 2.8915377

Fixed effects: dep ~ visit + pre + group
                Value Std.Error  DF   t-value p-value
(Intercept) 8.243741  3.247253 233  2.538682  0.0118
      visit -1.206358  0.167118 233 -7.218614  <.0001
        pre  0.468243  0.149474  58  3.132615  0.0027
      group -4.034921  1.121173  58 -3.598840  0.0007

 Correlation:
      (Intr)  visit    pre
visit -0.139
  pre -0.956  0.005
group -0.126 -0.047 -0.067

Standardized Within-Group Residuals:
       Min          Q1         Med         Q3      Max
 -3.315408 -0.5357005 -0.09072777 0.4617966 3.058502

Number of Observations: 295 Number of Groups: 61
```

# Non Normal Data

In this case, we cannot always specify a likelihood with an arbitrary structure. We can define random effect models by introducing a random intercept and slope into the linear predictor (generalized linear mixed models). These models can be difficult to estimate (GLLAMM6).
In the GEE approach, we can specify any covariance structure and link function without specifying the joint distribution of the the repeated observations.

REM and GEE lead to different interpretations of between subject effects. In the first case, a between subject effect stands for the difference between subjects conditional on the same random effect, while the parameters of GEE represent the average difference between subject.

# References

- Laird,Ware paper on REM, (Biometrics 1982)

- Zeger, Liang, Albert, on GEE (Biometrics, 1988)

- Horton, Lipsitz, on GEE software, (The American Statistician, 1998)