

A review of multilevel modelling in SPSS

*Alastair H Leyland
MRC Social and Public Health Sciences Unit
University of Glasgow
4 Lilybank Gardens
Glasgow G12 8RZ*

August 2004

1. Introduction to the software

1.1 Background

This review is based on SPSS version 12.0. The three SPSS commands of interest for multilevel modelling are all contained in the Advanced Models module, these being MIXED and VARCOMP. (An additional procedure GLM fits repeated measures models; however, random effects cannot be included in repeated measures designs in version 12.0.) The Advanced Models add capability to the SPSS Base system to conduct a range of additional analyses including generalised linear models and Cox regression; they complement the capabilities of the popular SPSS Base system. A major statistical package, SPSS is available in several languages. Most commands are available either through the graphical user interface or through the use of command syntax.

1.2 Software and hardware requirements

SPSS Advanced Models 12.0 requires an installation of SPSS Base 12.0. The system requirements for SPSS Base 12.0 are:

- Microsoft Windows 98, Me, NT 4.0, 2000 or XP
- Pentium®-class processor
- 200MB hard drive space
- 128MB RAM minimum
- SVGA monitor

SPSS Advanced Models 12.0 requires an additional 20MB of hard drive space, with the exact nature of other requirements varying according to the platform.

1.3 Data input/output

Data files may be opened using the Open Data command from the File menu. The range of file types supported makes SPSS compatible with many other packages; data files can be in any one of the following formats:

- SPSS, including SPSS for Windows, Macintosh, UNIX and SPSS portable
- SYSTAT
- Excel
- Lotus 1-2-3
- SYLK (Symbolic Link)
- dBASE
- a variety of SAS formats
- text (ASCII; fixed width or delimited)

By selecting the Open Database command from the File menu the user also has the facility to create a query to read in data from any database format for which they have a driver. Data can also be pasted from the clipboard.

The commands used to open data files are:

<i>Data type</i>	<i>Command</i>	<i>Subcommand</i>
SPSS	GET FILE = filename	
SPSS portable	IMPORT FILE = filename	
SYSTAT	GET TRANSLATE FILE = filename	/TYPE = SYS
Excel 5 or later	GET DATA	/FILE = filename /TYPE = XLS
Excel	GET TRANSLATE FILE = filename	/TYPE = XLS
Lotus 1-2-3 or Symphony	GET TRANSLATE FILE = filename	/TYPE = WK
SYLK format	GET TRANSLATE FILE = filename	/TYPE = SLK
dBASE	GET TRANSLATE FILE = filename	/TYPE = DBF
SAS	GET SAS DATA = filename	
ASCII	GET DATA	/FILE = filename /TYPE = TXT
ASCII (tab delimited)	GET TRANSLATE FILE = filename	/TYPE = TAB
Data accessed with ODBC driver	GET DATA	/FILE = filename /TYPE = ODBC

The data can be saved in a similar variety of formats by choosing Save As from the File menu. Alternatively, command syntax can be written using the SAVE, EXPORT and SAVE TRANSLATE commands.

1.4 Interface features

SPSS commands are written using a syntax language. All commands begin with a *keyword* which is the name of the command. Many commands also take *subcommands* and some may require additional specifications. For the purpose of this review all SPSS language is written in upper case, whilst user-defined variables are in lower case. Most command lines can be split into two or more lines at any point where a space could normally be inserted. Commands can be run from either batch (also known as production) or interactive modes; any command in interactive mode must finish with a full stop (period).

Although the use of syntax is essential to many users who want to ensure the replicability of their research, most SPSS commands are available through pointing and clicking in the menu-driven graphical user interface. The two commands identified as relating to multilevel modelling are available under the Analyze menu. MIXED can be found by selecting Mixed Models and then Linear (the only option available under Mixed Models in version 12.0). The VARCOMP command is obtained through Generalized Linear Model, selecting Variance Components.

2. Standard modelling tools for multilevel analysis

2.1 Fitting variance components using the VARCOMP command

The syntax of the VARCOMP command is

```
VARCOMP dependent variable BY factor list [WITH covariate list]
  /RANDOM = factor [factor...]
  [/METHOD = {MINQUE({1}) *}]
              {0}
              {ML }
              {REML }
              {SSTYPE({3}) }
              {1}
  [/INTERCEPT = {INCLUDE*}]
              {EXCLUDE }
  [/MISSING = {EXCLUDE*}]
              {INCLUDE }
  [/REGWGT = varname]
  [/CRITERIA = [CONVERGE({1.0E-8*})] [EPS({1.0E-8*})] [ITERATE({50*})]
              {n } {n } {n }
  [/PRINT = [EMS] [HISTORY({1*})] [SS]
              {n }]
```

```

[/OUTFILE = [VAREST] [{COVB}] (filename)
                {CORB}
[/DESIGN = {[INTERCEPT] [effect effect ...]]]

```

*Default if subcommand or keyword is omitted.

The VARCOMP command requires the higher level units to be specified as a factor in the main command line, with these units then specified as random effects (random factors) using the RANDOM subcommand. Further factors and covariates can be included in the main command. The default estimation method is the minimum norm quadratic unbiased estimator with unit prior weights. Alternative estimation methods – specified using the METHOD subcommand – are maximum likelihood (ML), restricted maximum likelihood (REML), or ANOVA method based on type I or type III sum of squares (SSTYPE(1) or SSTYPE(3)). The INTERCEPT subcommand determines whether or not the intercept is to be included in the model, and the MISSING subcommand determines the treatment of missing values. The REGWGT subcommand is used to specify regression weights in a weighted least squares regression model. CRITERIA is used to specify the convergence criterion in terms of the relative change in the objective function between iterations (CONVERGE), the tolerance for checking for singularity (EPS) and the maximum number of iterations (ITERATE). The PRINT subcommand is used to request output in terms of the objective function and variance components estimates at every n iterations (HISTORY(n), available only for maximum likelihood or restricted maximum likelihood estimation), the expected mean squares and the sums of squares (EMS and SS respectively, both available only for ANOVA estimation). OUTFILE is used to save the results of the estimation; VAREST will save the variance components estimates and COVB and CORB the covariance and correlation matrices (for ML and REML estimation only). The DESIGN subcommand is used to specify the effects (including interactions) included in a model, drawing from variables specified in the main command. The default is to include the intercept all covariates on the variable list, the main factorial effects and all orders of factor-by-factor interaction. Note that the VARCOMP procedure therefore provides only estimates of the variance components, not estimates of the regression coefficients. For this reason the rest of the review concentrates on the more general MIXED command.

2.2 Fitting multilevel models using the MIXED command

The syntax of the MIXED command is

```

MIXED dependent variable [BY factor list] [WITH covariate list]
  [/CRITERIA = [CIN({95*})] [MXITER({100*})] [MXSTEP({10*})] [SCORING({1*})]
                {n } {n } {n } {n }
                [SINGULAR({1E-12*})]
                {n }
                [{HCONVERGE({0*} {ABSOLUTE*}) }
                {n } {RELATIVE }
                {LCONVERGE({0*} {ABSOLUTE*}) }
                {n } {RELATIVE }
                {PCONVERGE({1E-6*} {ABSOLUTE*}) }
                {n } {RELATIVE }
[/EMMEANS = TABLES({OVERALL }
                    {factor }
                    {factor*factor...})
                [WITH(covariate={n } [covariate={n }...])
                    {MEAN} {MEAN}
                [COMPARE [(factor)] [REFCAT({n })] [ADJ({LSD* }
                    {FIRST} {BONFERRONI}
                    {LAST} {SIDAK }
[/FIXED = [effect [effect...]] [| [NOINT] [SSTYPE({1 })] ] ]
                {3*}
[/METHOD = {ML }
            {REML*}
[/MISSING = {EXCLUDE*}
            {INCLUDE }

```

```

[/PRINT = [CORB] [COVB] [CPS] [DESCRIPTIVES] [G] [HISTORY({1*})] [LMATRIX] [R]
                                     {n }
      [SOLUTION] [TESTCOV] ]
[/RANDOM = effect [effect...]
      [| [SUBJECT(varname[*varname[*...]])] [COVTYPE({VC*
                                     })]] ]
                                     {covstruct†}
[/REGWGT = varname]
[/REPEATED = varname[*varname[*...]] | SUBJECT(varname[*varname[*...]])
      [COVTYPE({DIAG*
      })]] ]
      {covstruct†}
[/SAVE = [tempvar [(name)] [tempvar [(name)]] ...] ]
[/TEST[(valuelist)]=['label'] effect valuelist ... [| effect valuelist ...] [divisor=n]]
      [; effect valuelist ... [| effect valuelist ...] [divisor=n] ]
[/TEST[(valuelist)] = ['label'] ALL list [| list] [divisor=n]
      [; ALL list [| list] [divisor=n]] ]

```

*Default if subcommand or keyword is omitted.

†covstruct can take one of the following values: AD1, AR1, ARH1, ARMA11, CS, CSH, CSR, DIAG, FA1, FAH1, HF, ID, TP, TPH, UN, UNR, VC.

The MIXED procedure can be used to fit a variety of mixed linear models including multilevel models. The command line is used to identify the dependent variable together with any factors and covariates to be included in the analysis. Note that, unlike the VARCOMP command, the MIXED command line does not require the specification of higher level units as factors. The CRITERIA subcommand is used to control the algorithm used for estimation and associated tolerance. Convergence can be determined by reference to the Hessian (HCONVERGE), the log-likelihood function (LCONVERGE) or the parameter estimates (PCONVERGE). The EMMEANS subcommand is used to provide the estimated marginal means for specific factors (or an overall mean if TABLES (OVERALL) is specified). The subcommand FIXED is used to specify which of the factors and covariates are to be included as fixed effects. Interactions can be included by using the BY keyword or, alternatively, an asterisk (*). An intercept is included unless the NOINT keyword is used. The METHOD subcommand is used to specify whether estimation is maximum likelihood or restricted maximum likelihood (the default), and the MISSING subcommand determines the treatment of missing values. The PRINT subcommand dictates the output of the MIXED analysis; options include printing the correlation and covariance matrices of the fixed parameter estimates (CORB and COVB), summary statistics of the dependent variable and any covariates for all combinations of factors including the higher level units specified using the RANDOM subcommand, the covariance matrix of the random effects (G), fixed and random parameter estimates (SOLUTION) and standard errors and Wald tests for the covariance parameters (TESTCOV). The RANDOM subcommand is used to specify the random part of the model; it specifies which factors or covariates are to be treated as random effects and at which level. To include a random intercept the keyword INTERCEPT must be specified as the first random effect in the RANDOM subcommand (the default is to exclude the intercept). There are a number of ways of specifying some models; for example, a random intercept model with higher level units defined by “L2_units” can be specified either by declaring “L2_units” to be a factor on the command line and entering “L2_units” as a random factor:

```

MIXED yvar BY L2_units
      /RANDOM = L2_units .

```

or by entering a random intercept and using the SUBJECT keyword of the RANDOM subcommand to identify the higher level units:

```

MIXED yvar
      /RANDOM = INTERCEPT | SUBJECT(L2_units) .

```

The RANDOM subcommand can be called repeatedly to configure complex random structures. The COVTYPE keyword specifies which of a list of pre-defined covariance structures is to be used. Many of the covariance structures allowed will be of interest for fitting growth curve or repeated measures models (e.g. first order ante-dependence AD1, first order autoregressive AR1, and diagonal or heterogeneous variances DIAG). For random effect models the common

choice will be an unstructured covariance matrix (UN) which will fit all variances and covariances between random effects. The REGWGT subcommand can be used to apply regression weights to the analysis. The REPEATED subcommand can be used to specify the covariance structure at level 1 in much the same way as the RANDOM subcommand. The SUBJECT keyword must be used to identify the hierarchical structure and must contain all of the variables specified as the subject (using the SUBJECT keyword) in any RANDOM subcommands. The SAVE subcommand can be used to save various case-specific statistics depending on the keywords used; the options are to save any of the predicted values based on the fixed part of the model (FIXPRED) or the fixed and higher level random parts (PRED), together with their standard errors (SEFIXP and SEPRED) and Satterthwaite degrees of freedom (DFFIXP and DFPRED) and the (composite level 1) residuals (RESID). Finally, the TEST subcommand allows the specification of null hypotheses as linear combination of parameters for both the fixed and random parts of the model. This subcommand conducts the F-test proposed by Fai and Cornelius (1996).

2.3 Information criteria available through the MIXED command

In addition to $-2 \log$ likelihood (or -2 restricted log likelihood if the METHOD is set to REML), the MIXED command gives four information criteria to assist model selection and comparison. These are the Akaike information criterion or AIC (Akaike, 1973):

$$AIC = -2l + 2d$$

the finite sample corrected AIC, or AICC (Hurvich and Tsai, 1989):

$$AIC_c = -2l + \frac{2dn}{(n-d-1)}$$

the consistent AIC (CAIC; Bozdogan, 1987):

$$CAIC = -2l + d[\log(n) + 1]$$

and the Bayesian information criterion or BIC (Schwarz, 1978):

$$BIC = -2l + d \log(n)$$

When using maximum likelihood estimation, n is taken to be the total number of level 1 units and d the number of fixed parameters plus the number of random parameters. For REML estimation, n is taken to be the total number of level 1 units minus the number of fixed parameters and d the number of random parameters.

2.4 Algorithm used by the MIXED command

The MIXED command uses Newton and scoring algorithms to maximise the likelihood (or restricted likelihood in the case of REML estimation). The algorithms used are those outlined by Wolfinger et al. (1994).

3. Model specifications – basic models

3.1 2-level normal models

The dataset used to illustrate the fitting of a 2-level normal model is the example taken from the user's guide to MLwiN (Rasbash et al., 2000), and comprises a sample of examination results from schools in six inner London Education Authorities (school boards). There are results for 4,059 students (level 1) nested within 65 schools (level 2), with between 2 and 198 students per school. The outcome for the i^{th} pupil in the j^{th} school, y_{ij} , has been standardised to a normal score with zero mean and unit variance, as has one of the covariates relating to prior ability (the London Reading Test score, x_{1ij}). The other individual covariate is the sex of the student, taking the value 0 for boys and 1 for girls. There is also an indicator of the sex mix in the school, taking the value 1 for a mixed school, 2 for a boys' school and 3 for a girls' school. Fitting such a categorical variable requires the creation of two dummy variables, x_{3j} and x_{4j} , indicating membership of two of these categories (the third category being used as the baseline or comparison group). These data have already been read into SPSS.

We start by fitting a variance components model to the data to estimate the effect of the covariates on the standardised exam score and to partition the variance between that arising due to differences between schools and that due to differences between students within schools. This model can be written as

$$\begin{aligned}y_{ij} &= \beta_0 + \beta_1 x_{1ij} + \beta_2 x_{2ij} + \beta_3 x_{3j} + \beta_4 x_{4j} + u_{0j} + e_{0ij} \\u_{0j} &\sim N(0, \sigma_{u0}^2) \\e_{0ij} &\sim N(0, \sigma_{e0}^2)\end{aligned}\tag{1}$$

The parameters that we will estimate are the 5 fixed parameters β_0, \dots, β_4 and the two variances σ_{u0}^2 and σ_{e0}^2 .

SPSS fits categorical variables as factors through the use of the BY keyword of the VARCOMP and MIXED commands. This creates the dummy variables necessary, using the last category as the comparison group. The reference category can therefore be changed by using the RECODE command. In equation (1) x_{2ij} has been coded as a dummy variable indicating the mean effect of girls relative to boys (so for SPSS we have $x_{2ij} = 1$ indicating a girl, $x_{2ij} = 2$ indicating a boy), and x_{3j} and x_{4j} indicate whether the school was a boys' school or a girls' school respectively. (Note that, in general, such factors can be numeric or string variables.)

We can fit model (1) using the code

```
MIXED normexam BY sex schlsex WITH standlrt
  /EMMEANS TABLES (sex*schlsex) WITH (standlrt=0)
  /FIXED = standlrt sex schlsex
  /RANDOM = INTERCEPT | SUBJECT (school)
  /METHOD = ML
  /PRINT = COVB G HISTORY SOLUTION TESTCOV
  /SAVE = FIXPRED (fix_pred) PRED (tot_pred) RESID (resid) .
EXECUTE .
```

The EMMEANS and PRINT subcommands are not required for this analysis – the code above is intended to illustrate their use. The SAVE command requests the predicted values from the fixed part of the model (saved in `fix_pred`), the predicted values from the fixed and random parts of the model (`tot_pred`) and the level 1 residuals (`resid`); this could again be omitted.

The estimates for this model (using both maximum likelihood ML and restricted maximum likelihood REML) estimation are given in table 1.

The next model includes an interaction between the two student level variables, the London Reading Test score and gender, in the fixed part of the model.

$$y_{ij} = \beta_0 + \beta_1 x_{1ij} + \beta_2 x_{2ij} + \beta_3 x_{3j} + \beta_4 x_{4j} + \beta_5 x_{1ij} x_{2ij} + u_{0j} + e_{0ij} \quad (2)$$

In equation (2) the parameter β_5 fits the difference between the slope with the London Reading Test score for girls (compared to boys). This model can be fitted by including an interaction term in the FIXED subcommand

```
/FIXED = standlrt sex standlrt*sex schlsex
```

The resulting parameter estimates are again shown in table 1.

The next model extends model (2) by allowing the coefficient of the London Reading Test score to vary at random across schools.

$$y_{ij} = \beta_0 + \beta_1 x_{1ij} + \beta_2 x_{2ij} + \beta_3 x_{3j} + \beta_4 x_{4j} + \beta_5 x_{1ij} x_{2ij} + u_{0j} + u_{1j} x_{1ij} + e_{0ij}$$

$$\begin{bmatrix} u_{0j} \\ u_{1j} \end{bmatrix} \sim N \left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} \sigma_{u0}^2 & \sigma_{u01} \\ \sigma_{u01} & \sigma_{u1}^2 \end{bmatrix} \right) \quad (3)$$

This model can be specified by changing the RANDOM subcommand

```
/RANDOM = INTERCEPT standlrt | SUBJECT(school) COVTYPE(UN)
```

Note that the default covariance structure for the RANDOM subcommand is variance components (VC); this means that, if the specification COVTYPE(UN) is omitted from the above command, SPSS will by default fit independent variances for the intercept and prior ability (i.e. the covariance term σ_{u01} in (3) will be omitted). The parameter estimates are given in table 1.

It does not appear possible in SPSS to model heterogeneity by fitting different variances for men and women:

$$y_{ij} = \beta_0 + \beta_1 x_{1ij} + \beta_2 x_{2ij} + \beta_3 x_{3j} + \beta_4 x_{4j} + \beta_5 x_{1ij} x_{2ij}$$

$$+ u_{0j} + u_{1j} x_{1ij} + e_{2ij} x_{2ij} + e_{6ij} x_{6ij} \quad (4)$$

$$\begin{bmatrix} e_{2ij} \\ e_{6ij} \end{bmatrix} \sim N \left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} \sigma_{e2}^2 & 0 \\ 0 & \sigma_{e6}^2 \end{bmatrix} \right)$$

where x_{6ij} is an indicator variable taking the value 1 for boys, 0 for girls i.e. $x_{6ij} = 1 - x_{2ij}$. Such a model can be fitted in other standard software packages such as MLwiN (Rasbash et al., 2000) and SAS (SAS Institute Inc., 1999).

3.2 3-level normal models

The dataset used to illustrate the 3-level normal response model is that previously analysed by Fielding et al.(2003) and refers to A/AS level examinations. The results for a Chemistry exam, in terms of the point score (0, 2, 4, ... 10) are given for 31,022 individuals from 2280 schools in 131 Local Education Authorities in England. The covariate we use for student i from school j in Education Authority k is an intake score (average GCSE score, x_{1ijk}). The model we consider is a variance components model:

$$\begin{aligned}
y_{ijk} &= \beta_0 + \beta_1 x_{ijk} + v_{0k} + u_{0jk} + e_{0ijk} \\
v_{0k} &\sim N(0, \sigma_{v_0}^2) \\
u_{0jk} &\sim N(0, \sigma_{u_0}^2) \\
e_{0ijk} &\sim N(0, \sigma_{e_0}^2)
\end{aligned}
\tag{5}$$

The addition of further levels to a model can be accomplished by using multiple RANDOM subcommands:

```

MIXED chem WITH gcse
  /FIXED = gcse
  /RANDOM = INTERCEPT | SUBJECT(lea) COVTYPE(ID)
  /RANDOM = INTERCEPT | SUBJECT(school*lea) COVTYPE(ID)
  /METHOD = ML
  /PRINT = COVB G HISTORY SOLUTION TESTCOV
  /SAVE = FIXPRED (fix_pred) PRED (tot_pred) RESID (resid) .
EXECUTE .

```

The declaration of SUBJECT(school*lea) following the previous SUBJECT(lea) indicates that schools are nested within Local Education Authorities. If the schools are given unique identifiers 1,...,2280 then the second RANDOM subcommand could be replaced by

```

/RANDOM = INTERCEPT | SUBJECT(school) COVTYPE(ID)

```

This second RANDOM subcommand then fits a cross-classification of schools by Local Education Authorities (see section 4.1), but the unique labelling of schools means that it fits a model identical to that specified in the previous syntax – the only difference being that the cross-classified model takes about 16 hours to converge instead of under 2 minutes for the nested model.

When fitting the nested model it makes no difference in theory whether unique identifiers are used for the schools or not; in practice, however, the use of unique identifiers for a large dataset such as this one is likely to result in memory problems. The solution is to recode the school identifier such that it runs from 1,..., n_k within each Local Education Authority. If the original coding of the schools is 1,...,2280 then one way of recoding them is as follows:

```

AGGREGATE OUTFILE = 'Temp.sav'
  /BREAK = lea
  /minschl = MIN(school) .
MATCH FILES /FILE = *
  /TABLE = 'Temp.sav'
  /BY lea .
COMPUTE school = school - minschl + 1 .
EXECUTE .
DELETE VARIABLES minschl .
EXECUTE .

```

3.3 Models for repeated measures data

Repeated measures models can be fitted using the MIXED command to balanced or unbalanced datasets, with or without time variant covariates. The REPEATED subcommand is used to specify the observations and the hierarchy (in addition to the RANDOM subcommand) as well as the covariance structure.

The data used to illustrate the repeated measures models is that analysed by Goldstein et al. (1994) and refer to the height of 26 boys aged 11 to 13 measured over 9 occasions

approximately 3 months apart. The data are balanced i.e. there are exactly 9 measurements made on each boy with no missing values.

We can first consider fitting a quartic polynomial to the height (cm) of the j^{th} boy measured on occasion i (at age t_{ij} , centred around 12 years), y_{ij} , with the coefficients of the intercept, linear and quadratic terms varying at random across the boys:

$$y_{ij} = \beta_0 + \beta_1 t_{ij} + \beta_2 t_{ij}^2 + \beta_3 t_{ij}^3 + \beta_4 t_{ij}^4 + u_{0j} + u_{1j} t_{ij} + u_{2j} t_{ij}^2 + e_{0ij}$$

$$\begin{bmatrix} u_{0j} \\ u_{1j} \\ u_{2j} \end{bmatrix} \sim N \left(\begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} \sigma_{u0}^2 & \sigma_{u01} & \sigma_{u02} \\ \sigma_{u01} & \sigma_{u1}^2 & \sigma_{u12} \\ \sigma_{u02} & \sigma_{u12} & \sigma_{u2}^2 \end{bmatrix} \right) \quad (6)$$

$$e_{0ij} \sim N(0, \sigma_{e0}^2)$$

This model can be fitted using:

```
MIXED height WITH age
  /FIXED = age age*age age*age*age age*age*age*age
  /RANDOM = INTERCEPT age age*age | SUBJECT(id) COVTYPE(UN)
  /METHOD = ML
  /PRINT = COVB G HISTORY SOLUTION TESTCOV
  /SAVE = FIXPRED (fix_pred) PRED (tot_pred) RESID (resid) .
EXECUTE .
```

We can now extend model (6) to fit first order autoregressive AR(1) errors at level 1. Goldstein et al. (1994) found evidence of seasonal effects on height; to counter this we include the sine and cosine of a seasonal (calendar year) time component T_{ij} in the fixed part of the model. Our model then becomes:

$$y_{ij} = \beta_0 + \sum_{h=1}^4 \beta_h t_{ij}^h + \beta_5 \sin(T_{ij}) + \beta_6 \cos(T_{ij}) + u_{0j} + u_{1j} t_{ij} + u_{2j} t_{ij}^2 + e_{0ij}$$

$$\text{Var}(e_{ij}) = \sigma_{e0}^2 \quad (7)$$

$$\text{Cov}(e_{ij}, e_{i'j}) = \rho^{|i-i'|} \sigma_{e0}^2 \quad \text{if } i \neq i'$$

Fitting this model requires, in addition to the declaration of the additional fixed parameters, a REPEATED subcommand specifying the measurement occasion i (coded 1 to 9 for each subject and called 'occasion') and an AR(1) covariance matrix at level 1.

```
COMPUTE sint = sin(season) .
COMPUTE cost = cos(season) .
MIXED height WITH age sint cost
  /CRITERIA = MXSTEP(25)
  /FIXED = age age*age age*age*age age*age*age*age sint cost
  /RANDOM = INTERCEPT age age*age | SUBJECT(id) COVTYPE(UN)
  /REPEATED = occasion | SUBJECT(id) COVTYPE(AR1)
  /METHOD = ML
  /PRINT = COVB G HISTORY SOLUTION TESTCOV
  /SAVE = FIXPRED (fix_pred) PRED (tot_pred) RESID (resid) .
EXECUTE .
```

If the data are unbalanced – if there aren't the same number of observations for each individual – SPSS is still able to fit the above repeated measures models.

4. Model specifications – more complex models

4.1 Cross-classified random effects models

As shown in section 3.2 above it is easy to fit cross-classified multilevel models using the MIXED command simply by adding another RANDOM subcommand and declaring an additional (non-nested) hierarchy. The data used to illustrate such a model relate to the exam scores of 3435 16 year old students in Fife, Scotland, with a view to disentangling the effects on educational attainment of the 148 primary schools and 19 secondary schools attended. The score of the i^{th} child who attended primary school j and secondary school k , y_{ijk} , is modelled in terms of the student's sex x_{1ijk} (taking the value 1 for girls, 0 for boys) only:

$$\begin{aligned}
 y_{ijk} &= \beta_0 + \beta_1 x_{1ijk} + v_{0k} + u_{0j} + e_{0ijk} \\
 v_{0k} &\sim N(0, \sigma_{v0}^2) \\
 u_{0j} &\sim N(0, \sigma_{u0}^2) \\
 e_{0ijk} &\sim N(0, \sigma_{e0}^2)
 \end{aligned} \tag{8}$$

This model can be fitted using:

```

MIXED attain BY sex
  /FIXED = sex
  /RANDOM = INTERCEPT | SUBJECT(sid)
  /RANDOM = INTERCEPT | SUBJECT(pid)
  /METHOD = ML
  /PRINT = COVB G HISTORY SOLUTION TESTCOV .
EXECUTE .

```

where sid and pid are the identifying codes for secondary school and primary school respectively. Note that the order of the RANDOM subcommands is not important. (The order of the subcommands is not important for any of these models fitted using the MIXED command.) For this model SPSS v12.0 would not let me save predicted values or residuals.

A model with the gender effect varying randomly across primary school:

$$\begin{aligned}
 y_{ijk} &= \beta_0 + \beta_1 x_{1ijk} + v_{0k} + u_{1j} x_{1ijk} + u_{2j} x_{2ijk} + e_{0ijk} \\
 v_{0k} &\sim N(0, \sigma_{v0}^2) \\
 \begin{bmatrix} u_{1j} \\ u_{2j} \end{bmatrix} &\sim N\left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} \sigma_{u1}^2 & \sigma_{u12} \\ \sigma_{u12} & \sigma_{u2}^2 \end{bmatrix}\right) \\
 e_{0ijk} &\sim N(0, \sigma_{e0}^2)
 \end{aligned} \tag{9}$$

where $x_{2ijk} = 1 - x_{1ijk}$ is a dummy variable taking the value 1 for boys and 0 for girls, can be fitted by changing the relevant RANDOM subcommand:

```

/RANDOM = sex | SUBJECT(pid) COVTYPE(UN)

```

Since sex has been declared as a factor on the MIXED command, the above RANDOM subcommand will allow both factors (boys and girls) to vary at random across primary schools and so the INTERCEPT should not be included. The covariance type needs to be specified as unstructured (UN) to estimate the covariance term σ_{u01} as described in section 3.1.

4.2 Multivariate normal response models

The multiple response model can be thought of as an extension of a repeated measures model – instead of a number of measurements of the same item made at different points in time we have measurements of a number of different items. We can use fixed effects to control for differences in the means between responses and random effects to model the different variances, but the real advantage of fitting multivariate response models is the ability to model the correlation between responses.

The data used to illustrate this model are examination scores for 1905 16 year old students from 73 schools in England, where results are available both for a written paper y_{Wjk} and for coursework y_{Cjk} for pupil j in school k . The fitted model is then

$$\begin{aligned}
 y_{Wjk} &= \beta_{W0} + \beta_{W1}x_{1jk} + v_{Wk} + u_{Wjk} \\
 y_{Cjk} &= \beta_{C0} + \beta_{C1}x_{1jk} + v_{Ck} + u_{Cjk} \\
 \begin{bmatrix} v_{Wk} \\ v_{Ck} \end{bmatrix} &\sim N\left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} \sigma_{Wv}^2 & \sigma_{WCv} \\ \sigma_{WCv} & \sigma_{Cv}^2 \end{bmatrix}\right) \\
 \begin{bmatrix} u_{Wjk} \\ u_{Cjk} \end{bmatrix} &\sim N\left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} \sigma_{Wu}^2 & \sigma_{WCu} \\ \sigma_{WCu} & \sigma_{Cu}^2 \end{bmatrix}\right)
 \end{aligned} \tag{10}$$

where x_{1jk} is a dummy variable taking the value 1 for boys, 0 for girls. The trick to fitting such a model is to stack the responses into a single column y_{ijk} and introduce an indicator variable I_{ijk} taking the value 1 for the written exam ($i = W$), 0 otherwise. Then we can write

$$y_{ijk} = I_{ijk}y_{Wjk} + (1 - I_{ijk})y_{Cjk} \tag{11}$$

If the data have multiple responses per record they need to be transformed from a format such as:

school	student	sex	writnexm	courswk
2	37	2	33	47.2
2	38	1	64	.

to a format with one response per record:

school	student	sex	index	y
2	37	2	1	33
2	37	2	2	47.2
2	38	1	1	64
2	38	1	2	.

This can be done with the following syntax:

```

COMPUTE index = 2 .
AGGREGATE OUTFILE = 'Temp.sav'
  /BREAK = school student
  /sex = MEAN(sex)
  /index = MEAN(index)
  /y = MEAN(courswk) .
COMPUTE y = writnexm .
COMPUTE index = 1 .
EXECUTE .
DELETE VARIABLES writnexm courswk .
EXECUTE .

```

```

ADD FILES /FILE = *
  /FILE = 'Temp.sav' .
SORT CASES BY school student (A) .
EXECUTE .

```

Note that the data may contain missing values; in the above example there is no score for coursework for student 38 in school 2. As mentioned in section 3.3 SPSS can analyse unbalanced repeated measures data, and since we use the REPEATED subcommand here this extends to missing multivariate responses.

To fit the model in SPSS we declare a 3-level model, with schools at the highest level and repeated measures on students at levels 1 and 2. There is, however, no modelling of the variance at the student level (there is no RANDOM subcommand with the keyword SUBJECT(student)).

```

MIXED y BY index sex
  /FIXED = index index*sex | NOINT
  /RANDOM = index | SUBJECT(school) COVTYPE (UN)
  /REPEATED = index | SUBJECT(school*student) COVTYPE (UN)
  /METHOD = ML
  /PRINT = COVB G HISTORY SOLUTION TESTCOV
  /SAVE = FIXPRED (fix_pred) PRED (tot_pred) RESID (resid) .
EXECUTE .

```

The use of school*student on the REPEATED subcommand indicates the nesting of students within schools as discussed in section 3.2, and assumes that the students are numbered from 1 to n_k within each school. Specifying no intercept on the FIXED subcommand (using the NOINT keyword) ensures that separate estimates are obtained for each of the written exam and coursework components. Similarly there is no intercept included on the RANDOM subcommands – the variable index distinguishes between the two responses and will therefore fit separate random intercepts. Using the COVTYPE (UN) keyword specifies an unstructured covariance matrix for responses and schools.

This model can be extended and generalised as required.

5. Obtaining residuals in SPSS

SPSS 12.0 does not provide the higher level residuals directly, presumably because these are seen as some kind of nuisance terms. However, in many cases there will be substantive interest in the residuals and the SAVE subcommand can be used to save the fixed part predictions (FIXPRED) as well as the predictions from the fixed and random parts of the model (PRED). For the general linear multilevel model, written in matrix form,

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\boldsymbol{\gamma} + \mathbf{e} \quad (12)$$

where $\boldsymbol{\gamma}$ is a stacked vector of all residuals (slopes and intercepts) at all levels and \mathbf{Z} is the corresponding design matrix, the predictions from the fixed part correspond to

$$\hat{\mathbf{Y}}^* = \mathbf{X}\hat{\boldsymbol{\beta}} \quad (13)$$

and the predictions from the fixed and random parts are given by

$$\hat{\mathbf{Y}} = \mathbf{X}\hat{\boldsymbol{\beta}} + \mathbf{Z}\hat{\boldsymbol{\gamma}} \quad (14)$$

It follows from (13) and (14) that the predicted residuals $\hat{\boldsymbol{\gamma}}$ are given by

$$\hat{\gamma} = (\mathbf{Z}^T \mathbf{Z})^{-1} \mathbf{Z}^T (\hat{\mathbf{Y}} - \hat{\mathbf{Y}}^*) \quad (15)$$

For the trivial example of a 2-level variance components model given by (1) or (2), the following code uses the MATRIX command (and, in particular, the SOLVE function) in SPSS to obtain estimates of the school-level residuals.

```
AUTORECODE VARIABLES = school
/INTO l2id .
SORT CASES BY l2id .
* get composite residuals .
COMPUTE comp_res = tot_pred - fix_pred .
* make sure MXLOOP is greater than the number of schools .
SET MXLOOP = 100 .
MATRIX .
  GET l2id
  /FILE = *
  /VARIABLES = l2id .
  GET school
  /FILE = *
  /VARIABLES = school .
  GET comp_res
  /FILE = *
  /VARIABLES = comp_res .
  COMPUTE temp_mat = (l2id = 1) .
  COMPUTE zmat = {temp_mat} .
  LOOP i = 2 TO l2id(NROW(l2id)) .
    COMPUTE temp_mat = (l2id = i) .
    COMPUTE zmat = {zmat, temp_mat} .
  END LOOP .
  COMPUTE zTz = T(zmat)*zmat .
  COMPUTE zTy = T(zmat)*comp_res .
  COMPUTE res_2 = SOLVE(zTz,zTy) .
  COMPUTE zTy = T(zmat)*school .
  COMPUTE schl_2 = SOLVE(zTz,zTy) .
  SAVE {schl_2,res_2}
  /OUTFILE = *
  /VARIABLES = school res_2_1 .
END MATRIX .
EXECUTE .
```

The MATRIX command of this code can be modified to estimate, for example, residuals for the 2-level random slopes model given by (3):

```
MATRIX .
  GET l2id
  /FILE = *
  /VARIABLES = l2id .
  GET school
  /FILE = *
  /VARIABLES = school .
  GET comp_res
  /FILE = *
  /VARIABLES = comp_res .
  GET standlrt
  /FILE = *
  /VARIABLES = standlrt .
  COMPUTE temp_mat = (l2id = 1) .
  COMPUTE zmat = {temp_mat} .
  LOOP i = 2 TO l2id(NROW(l2id)) .
    COMPUTE temp_mat = (l2id = i) .
```

```

    COMPUTE zmat = {zmat, temp_mat} .
END LOOP .
COMPUTE zTz = T(zmat)*zmat .
COMPUTE zTy = T(zmat)*school .
COMPUTE schl_2 = SOLVE(zTz,zTy) .
LOOP i = 1 TO l2id(NROW(l2id)) .
    COMPUTE temp_mat = (l2id = i)&*standlrt .
    COMPUTE zmat = {zmat, temp_mat} .
END LOOP .
COMPUTE zTz = T(zmat)*zmat .
COMPUTE zTy = T(zmat)*comp_res .
COMPUTE res_2 = SOLVE(zTz,zTy) .
COMPUTE temp_mat = IDENT(l2id(NROW(l2id)),2*l2id(NROW(l2id))) .
COMPUTE res_2_1 = temp_mat*res_2 .
COMPUTE temp_mat = {0*IDENT(l2id(NROW(l2id))),
IDENT(l2id(NROW(l2id)))} .
COMPUTE res_2_2 = temp_mat*res_2 .
SAVE {schl_2,res_2_1,res_2_2}
/OUTFILE = *
/VARIABLES = school res_2_1 res_2_2 .
END MATRIX .
EXECUTE .

```

Of course the residuals are of little use in themselves without their corresponding standard errors. The dispersion matrix of the residuals can be estimated using formulae given by e.g. Goldstein (2003).

6. Conclusions

Multilevel modelling in SPSS has definite limitations; in particular, the restriction to normal response models means that several classes of model cannot be fitted. These include such common models as multilevel logistic regression and multilevel Poisson regression models and, through these, developments such as multilevel categorical responses or multilevel Cox regression.

The major limitation to the normal response models is the restricted ability to specify the covariance matrix at the lowest level. In particular, this means that SPSS is not able to fit models with heterogeneous variances as in equation (4). This may seem like a minor limitation but in effect it means that the user must hypothesise that the lowest level variance is the same for all subgroups (and that it is independent of the value of any covariate) without being able to test these hypotheses. This becomes particularly important when testing for random slopes at higher levels, since the inability to model the variance at the lowest level may effect the outcome of such tests. Moreover, although it is possible to obtain higher level residuals from the models that SPSS fits, it is unduly cumbersome at present.

However, there are some strengths to the SPSS MIXED command. The alteration or addition of RANDOM subcommands makes it easy to change the random specification of a model (at the higher levels) or to add further levels, and it is as straightforward to fit cross-classified models as it is to fit hierarchical models. The REPEATED subcommand provides a wide range of correlation functions, and the use of these makes it simple to fit normal multivariate response models. There is no requirement for datasets to be balanced or complete, the information criteria provided are fairly comprehensive and the algorithm used is fast. The MIXED command is also available through the Windows interface (as opposed to through the use of the command syntax); a description of the use of the MIXED command through the Windows interface can be found elsewhere (Landau and Everitt, 2004).

The widespread use of SPSS means that, if it to be taken seriously as a statistical package, it is important that multilevel data analysis should be available. The MIXED command already covers most of the multilevel analyses that most users will require for (normally distributed)

continuous outcomes. However, in many disciplines continuous measures will be the exception rather than the rule and SPSS will remain limited until it introduces commands to fit generalised discrete response multilevel models. Put it like this: unless all of your (multilevel) data have normally distributed responses you are going to need to use a package other than SPSS to analyse them. In which case, is it worth taking the time to learn how to use the MIXED command in SPSS when you are also going to have to learn to use other software?

References

- Akaike H. (1973) Information theory and an extension of the maximum likelihood principle. In: Petrov BN, Csaki F, eds. *2nd International Symposium on Information Theory*. Budapest: Akademiai Kiado, 267-281.
- Bozdogan H. (1987) Model selection and Akaike's information criterion (AIC): the general theory and its analytical extensions. *Psychometrika* **52**, 345-370.
- Fai A. H. T., Cornelius P. L. (1996) Approximate F-tests of multiple degree of freedom hypotheses in generalized linear least squares analyses of unbalanced split-plot experiments. *Journal of Statistical Computation and Simulation* **54**, 363-378.
- Fielding A., Yang M., Goldstein H. (2003) Multilevel ordinal models for examination grades. *Statistical Modelling* **3**.
- Goldstein H. (2003) *Multilevel statistical models*. London: Arnold.
- Goldstein H., Healy M. J. R., Rasbash J. (1994) Multilevel time series models with applications to repeated measures data. *Statistics in Medicine* **13**, 1643-1655.
- Hurvich C. M., Tsai C.-L. (1989) Regression and time series model selection in small samples. *Biometrika* **76**, 297-307.
- Landau S., Everitt B. S. (2004) *A Handbook of Statistical Analyses using SPSS*. Boca Raton: Chapman & Hall.
- Rasbash J., Browne W., Goldstein H., Yang M., Plewis I., Healy M., Woodhouse G., Draper D., Langford I., Lewis T. (2000) *A User's Guide to MLwiN*. London: Multilevel Models Project, Institute of Education, University of London.
- SAS Institute Inc. (1999) *SAS/STAT User's Guide, Version 7-1*. Cary, NC: SAS Institute Inc.
- Schwarz G. (1978) Estimating the dimension of a model. *Annals of Statistics* **6**, 461-464.
- Wolfinger R., Tobias R., Sall J. (1994) Computing Gaussian likelihoods and their derivatives for general linear mixed models. *SIAM Journal on Scientific Computing* **15**, 1294-1310.

Table 1: parameter estimates for 2-level models

Model	Parameter	ML			REML		
		Estimate	SE	Time	Estimate	SE	Time
(1)	β_0	-0.0091	0.0763	1s	-0.0094	0.0779	1s
	β_1	0.5600	0.0124		0.5598	0.0125	
	β_2	0.1672	0.0341		0.1674	0.0341	
	β_3	-0.1590	0.0873		-0.1590	0.0894	
	β_4	0.0187	0.1232		0.0187	0.1261	
	σ_{u0}^2	0.0811	0.0165		0.0858	0.0178	
	σ_{e0}^2	0.5623	0.0126		0.5625	0.0126	
	-2 log like	9325.43			9347.67		
(2)	β_0	-0.0091	0.0763	1s	-0.0094	0.0779	1s
	β_1	0.5628	0.0184		0.5626	0.0184	
	β_2	0.1672	0.0341		0.1673	0.0341	
	β_3	-0.1588	0.0873		-0.1588	0.0894	
	β_4	0.0188	0.1232		0.0189	0.1261	
	β_5	-0.0051	0.0246		-0.0051	0.0246	
	σ_{u0}^2	0.0811	0.0166		0.0859	0.0178	
	σ_{e0}^2	0.5623	0.0126		0.5627	0.0126	
-2 log like	9325.39			9353.20			
(3)	β_0	-0.0114	0.0728	2s	-0.0120	0.0742	2s
	β_1	0.5507	0.0255		0.5503	0.0257	
	β_2	0.1683	0.0338		0.1686	0.0338	
	β_3	-0.1784	0.0801		-0.1779	0.0821	
	β_4	-0.0007	0.1136		-0.0004	0.1163	
	β_5	0.0069	0.0294		0.0069	0.0295	
	σ_{u0}^2	0.0795	0.0164		0.0837	0.0175	
	σ_{u01}	0.0202	0.0067		0.0205	0.0070	
	σ_{u1}^2	0.0147	0.0046		0.0152	0.0047	
	σ_{e0}^2	0.5502	0.0124		0.5504	0.0124	
	-2 log like	9281.07			9308.24		

Table 2: parameter estimates for 3-level model

Model	Parameter	ML			REML		
		Estimate	SE	Time	Estimate	SE	Time
(5)	β_0	-9.9067	0.1089	112s	-9.9063	0.1090	120s
	β_1	2.4726	0.0169		2.4726	0.0169	
	σ_{v0}^2	0.0136	0.0135		0.0148	0.0139	
	σ_{u0}^2	1.1662	0.0555		1.1662	0.0555	
	σ_{e0}^2	5.1541	0.0431		5.1542	0.0555	
	-2 log like	141685.6			141728.0		

Table 3: parameter estimates for repeated measures models

Model	Parameter	ML			REML		
		Estimate	SE	Time	Estimate	SE	Time
(6)	β_0	148.9753	1.5396	1s	148.9753	1.5701	1s
	β_1	6.1659	0.3510		6.1658	0.3574	
	β_2	1.0906	0.3490		1.0905	0.3525	
	β_3	0.4678	0.1625		0.4678	0.1635	
	β_4	-0.3404	0.3002		-0.3404	0.3021	
	σ_{u0}^2	61.5486	17.0858		64.0120	18.1211	
	σ_{u1}^2	2.7627	0.7823		2.8748	0.8297	
	σ_{u2}^2	0.6304	0.2248		0.6604	0.2384	
	σ_{u01}	7.9922	3.0232		8.3119	3.2063	
	σ_{u02}	1.3633	1.4058		1.4158	1.4910	
	σ_{u12}	0.8747	0.3419		0.9096	0.3626	
	σ_{e0}^2	0.2175	0.0246		0.2203	0.0251	
	-2 log like	627.278			629.825		
	(7)	β_0	148.8878	1.5350	4s	148.8877	1.5656
β_1		6.2528	0.3741		6.2522	0.3812	
β_2		1.9434	0.5035		1.9437	0.5107	
β_3		0.1855	0.2410		0.1857	0.2444	
β_4		-1.1160	0.4307		-1.1162	0.4365	
β_5		-0.2736	0.0972		-0.2737	0.0986	
β_6		0.1212	0.0613		0.1212	0.0622	
σ_{u0}^2		60.6582	16.9393		63.0951	17.9734	
σ_{u1}^2		2.2110	0.6972		2.3051	0.7425	
σ_{u2}^2		0.3116	0.2730		0.3235	0.2861	
σ_{u01}		7.7376	2.8631		8.0443	3.0412	
σ_{u02}		2.0516	1.5335		2.1305	1.6179	
σ_{u12}		0.8061	0.3424		0.8390	0.3623	
σ_{e0}^2		0.7334	0.0837		0.7626	0.0882	
ρ		0.7084			0.7116		
-2 log like	622.136			631.602			

Table 4: parameter estimates for cross-classified model

Model	Parameter	ML			REML		
		Estimate	SE	Time	Estimate	SE	Time
(8)	β_0	5.2574	0.1807	12s	5.2552	0.1843	12s
	β_1	0.4986	0.0982		0.4985	0.0983	
	σ_{v0}^2	0.3457	0.1609		0.3697	0.1733	
	σ_{u0}^2	1.1043	0.2023		1.1096	0.2036	
	σ_{e0}^2	8.0534	0.1990		8.0551	0.1991	
	-2 log like	17123.5			17127.9		
	β_0	5.2605	0.1783	1m30s	5.2580	0.1821	1m37s
	β_1	0.4939	0.1072		0.4940	0.1078	
	σ_{v0}^2	0.3409	0.1602		0.3652	0.1727	
	σ_{u1}^2	1.2960	0.2810		1.3066	0.2837	
	σ_{u12}^2	1.0652	0.2085		1.0667	0.2100	
	σ_{u2}^2	1.0258	0.2224		1.0324	0.2242	
	σ_{e0}^2	8.0062	0.2013		8.0050	0.2013	
	-2 log like	17121.4			17125.7		

Table 5: parameter estimates for multivariate response model

Model	Parameter	ML			REML		
		Estimate	SE	Time	Estimate	SE	Time
(8)	β_{W0}	49.0084	0.9318	3s	49.0096	0.9380	3s
	β_{C0}	69.6230	1.1719		69.6211	1.1795	
	β_{W1}	-2.4930	0.5603		-2.4913	0.5605	
	β_{C1}	6.7567	0.6706		6.7574	0.6709	
	σ_{Wv}^2	46.5648	9.3531		47.3794	9.5623	
	σ_{WCv}	24.9371	8.9916		25.3663	9.1903	
	σ_{Cv}^2	75.1936	14.6729		76.4476	14.9919	
	σ_{Wu}^2	124.4335	4.3363		124.5024	4.3400	
	σ_{WCu}	72.7489	4.1521		72.7841	4.1555	
	σ_{Cu}^2	180.0697	6.2499		180.1729	6.2553	
	-2 log like	26799.5			26794.6		