

*Statistics in Medicine(2002) . To appear.*

## **Multilevel Modelling of medical data**

By  
Harvey Goldstein  
William Browne  
And  
Jon Rasbash

Institute of Education, University of London<sup>1</sup>

### **Summary**

This tutorial presents an overview of multilevel or hierarchical data modelling and its applications in medicine. A description of the basic model for nested data is given and it is shown how this can be extended to fit flexible models for repeated measures data and more complex structures involving cross classifications and multiple membership patterns within the software package *MLwiN*. A variety of response types are covered and both frequentist and Bayesian estimation methods are described.

### **Keywords**

Complex data structures, mixed model, multilevel model, random effects model, repeated measures.

---

<sup>1</sup> Correspondence to Professor H. Goldstein, Institute of Education, 20 Bedford Way, London, WC1H 0AL, UK. Email: [h.goldstein@ioe.ac.uk](mailto:h.goldstein@ioe.ac.uk). FAX: +44 20 7612 6686. TEL: +44 20 7612 6652

## 1. Scope of tutorial

The tutorial covers the following topics

1. The nature of multilevel models with examples.
2. Formal model specification for the basic Normal (nested structure) linear multilevel model with an example.
3. The *MLwiN* software.
4. More complex data structures: Complex variance, multivariate models and cross-classified and multiple membership models.
5. Discrete response models, including Poisson, binomial and multinomial error distributions.
6. Specific application areas including survival models, repeated measures models, spatial models and meta analysis.
7. Estimation methods, including maximum and quasi likelihood, and MCMC.

Further information about multilevel modelling and software details can be obtained from the web site of the Multilevel Models Project, <http://multilevel.ioe.ac.uk/>

## 2. The nature of multilevel models

Traditional statistical models were developed making certain assumptions about the nature of the dependency structure among the observed responses. Thus, in the simple regression model  $y_i = \beta_0 + \beta_1 x_i + e_i$  the standard assumption is that the  $y_i$  given  $x_i$  are independently identically distributed (iid), and the same assumption holds also for generalised linear models. In many real life situations, however, we have data structures, whether observed or by design, for which this assumption does not hold.

Suppose, for example, that the response variable is the birthweight of a baby and the predictor is, say, maternal age, and data are collected from a large number of maternity units located in different physical and social environments. We would expect that the maternity units would have different mean birth weights, so that knowledge of the maternity unit already conveys some information about the baby. A more suitable model for these data is now

$$y_{ij} = \beta_0 + \beta_1 x_{ij} + u_j + e_{ij} \quad (1)$$

where we have added another subscript to identify the maternity unit and included a unit-specific effect  $u_j$  to account for mean differences amongst units. If we assume that the maternity units are randomly sampled from a population of units, then the unit specific effect is a random variable and (1) becomes a simple example of a 2-level model. Its complete specification, assuming Normality, can be written as follows

$$\begin{aligned} y_{ij} &= \beta_0 + \beta_1 x_{ij} + u_j + e_{ij} \\ u_j &\sim N(0, \sigma_u^2), \quad e_{ij} \sim N(0, \sigma_e^2) \\ \text{cov}(u_j, e_{ij}) &= 0 \\ \text{cov}(y_{i_1 j}, y_{i_2 j} \mid x_{ij}) &= \sigma_u^2 \geq 0 \end{aligned} \quad (2)$$

where  $i_1, i_2$  are two births in the same unit  $j$  with, in general, a positive covariance between the responses. This lack of independence, arising from two sources of variation at different levels of the data hierarchy (births and maternity units) contradicts the traditional linear model assumption and leads us to consider a new class of models. Model (2) can be elaborated in a number of directions, including the addition of further covariates or levels of nesting. An important direction is where the coefficient (and any further coefficients) is allowed to have a random distribution. Thus, for example the age relationship may vary across clinics and, with a slight generalisation of notation, we may now write (2) as

$$\begin{aligned} y_{ij} &= \beta_{0ij} x_{0ij} + \beta_{1j} x_{1ij} \\ \beta_{0ij} &= \beta_0 + u_{0j} + e_{0ij} \\ \beta_{1j} &= \beta_1 + u_{1j} \\ x_{0ij} &= 1 \\ \text{var}(u_{0j}) &= \sigma_{u0}^2, \quad \text{var}(u_{1j}) = \sigma_{u1}^2, \\ \text{cov}(u_{0j}, u_{1j}) &= \sigma_{u01}, \quad \text{var}(e_{0ij}) = \sigma_{e0}^2 \end{aligned} \quad (3)$$

and in later sections we shall introduce further elaborations. The regression coefficients  $\beta_0, \beta_1$  are usually referred to as ‘fixed parameters’ of the model and the set of variances and covariances as the random parameters. Model (3) is often referred to as a ‘random coefficient’ or ‘mixed’ model.

At this point we note that we can introduce prior distributions for the parameters of (3), so allowing Bayesian models. We leave this topic, however, for a later section where we discuss MCMC estimation.

Another, instructive, example of a 2-level data structure for which a multilevel model provides a powerful tool, is that of repeated measures data. If we measure the weight of a sample of babies after birth at successive times then the repeated occasion of measurement becomes the lowest level unit of a 2-level hierarchy where the individual baby is the level 2 unit. In this case model (3) would provide a simple description with  $x_{1ij}$  being time or age. In practice linear growth will be an inadequate description and we would wish to fit at least a (spline) polynomial function, or perhaps a non-linear function where several coefficients varied randomly across individual babies, that is each baby has its own growth pattern. We shall return to this example in more detail later, but for now note that an important feature of such a characterisation is that it makes no particular requirements for every baby to be measured at the same time points or for the time points to be equally spaced.

The development of techniques for specifying and fitting multilevel models since the mid 1980s has produced a very large class of useful models. These include models with discrete responses, multivariate models, survival models, time series models etc. In this tutorial we cannot cover the full range but will give references to existing and ongoing work that readers may find helpful. In addition the introductory book by Snijders and Bosker [1] and the edited collection of health applications by Leyland and Goldstein [2] may be found useful by readers.

A detailed introduction to the 2-level model with worked examples and discussion of hypothesis tests and basic estimation techniques is given in an earlier tutorial [3] that also gives details of two computer packages, HLM and SAS that can perform some of the analyses we describe in the present tutorial. The *MLwiN* software has been specifically developed for fitting very large and complex models, using both frequentist and Bayesian estimation and it is this particular set of features that we shall concentrate on.

### **3. Marginal vs hierarchical models**

At this stage it is worth emphasising the distinction between multilevel models and so called ‘marginal’ models such as the GEE model [4,5]. When dealing with hierarchical data these latter models typically start with a formulation for the covariance structure, for example but not necessarily based upon a multilevel structure such as (3), and aim to provide estimates with acceptable properties only for the fixed parameters in the model, treating the existence of any random parameters as a necessary ‘nuisance’ and without providing explicit estimates for them. More specifically, the estimation procedures used in marginal models are known to have useful asymptotic properties in the case where the exact form of the random structure is unknown.

If interest lies only in the fixed parameters, marginal models may be useful. Even here, however, they may be inefficient if they utilise a covariance structure that is substantially incorrect. They are, however, generally more robust than multilevel models to serious misspecification of the covariance structure [6]. Fundamentally, however, marginal models address different research questions. From a multilevel perspective, the failure explicitly to model the covariance structure of complex data is to ignore information about variability that, potentially, is as important as knowledge of the average or fixed effects. Thus, in the simple repeated measures example of baby weights knowledge of how individual growth rates vary between babies, possibly differentially according to say demographic factors, will be important data and in a later section we will show how such information can be used to provide efficient predictions in the case of human growth.

When we discuss discrete response multilevel models we will show how to obtain information equivalent to that obtained from marginal models. Apart from that the remainder of this paper will be concerned with multilevel models. For a further discussion of the limitations of marginal models see the paper by Lindsey and Lambert [7].

## 4. Estimation for the multivariate Normal model

We write the general Normal 2-level model as follows, with natural extensions to three or more levels:

$$\begin{aligned}
\mathbf{Y} &= \mathbf{X}\boldsymbol{\beta} + \mathbf{E} \\
\mathbf{Y} &= \{y_{ij}\}, \quad \mathbf{X} = \{\mathbf{X}_{ij}\}, \quad \mathbf{X}_{ij} = \{x_{0ij}, x_{1ij}, \dots, x_{pij}\} \\
\mathbf{E} &= \mathbf{E}^{(2)} + \mathbf{E}^{(1)} \\
\mathbf{E}^{(2)} &= \{\mathbf{E}_j^{(2)}\}, \quad \mathbf{E}_j^{(2)} = \mathbf{z}_j^{(2)} \mathbf{e}_j^{(2)}, \quad \mathbf{z}_j^{(2)} = \{z_{ij}^{(2)}\} \\
\mathbf{z}_{ij}^{(2)} &= \{z_{0j}^{(2)}, z_{1j}^{(2)}, \dots, z_{q_2j}^{(2)}\}, \quad \mathbf{e}_j^{(2)} = \{e_{0j}^{(2)}, e_{1j}^{(2)}, \dots, e_{q_2j}^{(2)}\}^T \\
\mathbf{E}^{(1)} &= \{\mathbf{E}_{ij}^{(1)}\}, \quad \mathbf{E}_{ij}^{(1)} = \mathbf{z}_{ij}^{(1)} \mathbf{e}_{ij}^{(1)} \\
\mathbf{z}_{ij}^{(1)} &= \{z_{0j}^{(1)}, z_{1j}^{(1)}, \dots, z_{q_2j}^{(1)}\}, \quad \mathbf{e}_{ij}^{(1)} = \{e_{0ij}^{(1)}, e_{1ij}^{(1)}, \dots, e_{q_1ij}^{(1)}\}^T \\
\mathbf{e}_j^{(2)} &= \{\mathbf{e}_j^{(2)}\}, \quad \mathbf{e}_{ij}^{(1)} = \{\mathbf{e}_{ij}^{(1)}\} \\
\mathbf{e}_j^{(2)} &\sim \mathbf{N}(\mathbf{0}, \boldsymbol{\Omega}_2), \quad \mathbf{e}_{ij}^{(1)} \sim \mathbf{N}(\mathbf{0}, \boldsymbol{\Omega}_{1j}) \\
&[\text{Typically } \boldsymbol{\Omega}_{1j} = \boldsymbol{\Omega}_1] \\
E(e_{hj}^{(2)} e_{h'j'}^{(2)})_{j \neq j'} &= E(e_{hj}^{(1)} e_{h'j'}^{(1)})_{i \neq i'} = E(e_{hj}^{(2)} e_{h'j'}^{(1)}) = 0 \\
&\text{yields the block diagonal structure} \\
\mathbf{V} &= E(\tilde{\mathbf{Y}}\tilde{\mathbf{Y}}^T) = \bigoplus_j (\mathbf{V}_{2j} + \mathbf{V}_{1j}) \\
\tilde{\mathbf{Y}} &= \mathbf{Y} - \mathbf{X}\boldsymbol{\beta} \\
\mathbf{V}_{2j} &= \mathbf{z}_j^{(2)} \boldsymbol{\Omega}_2 \mathbf{z}_j^{(2)T}, \quad \mathbf{V}_{1j} = \bigoplus_i \mathbf{z}_{ij}^{(1)} \boldsymbol{\Omega}_{1j} \mathbf{z}_{ij}^{(1)T}
\end{aligned} \tag{4}$$

In this formulation we allow any number of random effects or coefficients at each level; we shall discuss the interpretation of multiple level 1 random coefficients in a later section.

A number of efficient algorithms are available for obtaining maximum likelihood (ML) estimates for (4). One [8] is an iterative generalised least squares procedure (IGLS) that will also produce restricted maximum likelihood estimates (RIGLS or REML) and is formally equivalent to a Fisher scoring algorithm [9]. Note that RIGLS or REML should be used in small samples to correct for the underestimation of IGLS variance estimates. The EM algorithm can also be used [10,11]. Our examples use RIGLS (REML) estimates as implemented in the *MlwiN* software package [12] and we will also discuss Bayesian models. A simple description of the IGLS algorithm is as follows.

From (4) we have

$$\mathbf{V} = E(\tilde{\mathbf{Y}}\tilde{\mathbf{Y}}^T) = \oplus_j (\mathbf{V}_{2j} + \mathbf{V}_{1j})$$

$$\tilde{\mathbf{Y}} = \mathbf{Y} - \mathbf{X}\boldsymbol{\beta}$$

The IGLS algorithm proceeds by first carrying out a GLS estimation for the fixed parameters ( $\boldsymbol{\beta}$ ) using a working estimator of  $\mathbf{V}$ . The vectorised cross product matrix of ‘raw’ residuals  $\hat{\tilde{\mathbf{Y}}}\hat{\tilde{\mathbf{Y}}}^T$  where  $\hat{\tilde{\mathbf{Y}}} = \mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}$ , is then used as the response in a GLS estimation where the explanatory variable design matrix is determined by the last line of (4). This provides updated estimates for the  $\Omega_{1j}$  and  $\Omega_2$  and hence  $V$ . The procedure is repeated until convergence. In the simple case we have been considering so far where the level 1 residuals are iid, for a level 2 unit (individual) with just 3 level 1 units (occasions) there are just 6 distinct raw residual terms and the level 1 component  $\mathbf{V}_{1j}$  is simply  $\sigma_e^2 \mathbf{I}_3$ . Written as a vector of the lower triangle this becomes

$$\sigma_e^2 \begin{pmatrix} 1 \\ 0 \\ 1 \\ 0 \\ 0 \\ 1 \end{pmatrix} \quad (5)$$

and the vector of ones and zeroes becomes the level 1 explanatory variable for the GLS estimation, in this case providing the coefficient that is the estimator of  $\sigma_e^2$ . Similarly, for a model where there is a single variance term at level 2, the level 2 component  $\mathbf{V}_{2j}$  written as a lower triangle vector is

$$\sigma_u^2 \begin{pmatrix} 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \end{pmatrix}$$

Goldstein [13] shows that this procedure produces maximum likelihood estimates under Normality.

#### **4. The *MLwiN* software**

*MLwiN* has been under development since the late 1980s, first as a command-driven DOS based program, *MLn*, and since 1998 in a fully-fledged windows version, currently in release 1.10. It is produced by the Multilevel Models project based within the Institute of Education, University of London, and supported largely by project funds from the UK Economic and Social research Council. The software has been developed alongside advances in methodology and with the preparation of manuals and other training materials.

Procedures for fitting multilevel models are now available in several major software packages such as STATA, SAS and S plus. In addition there are some special purpose packages, which are tailored to particular kinds of data or models. MIXOR provides ML estimation for multicategory responses and HLM is used widely for educational data. See Zhou et al. [14] for a recent review and Sullivan et al. [3] for a description of the use of HLM and SAS. Many of the models discussed here can also be fitted readily in the general purpose MCMC software package BUGS [15].

*MLwiN* has some particular advanced features that are not available in other packages and it also has a user interface designed for fully interactive use. In later sections we will illustrate some of the special features and models available in *MLwiN* but first give a simple illustration of the user interface. We shall assume that the user wishes to fit the simple 2-level model given by (1).

In this tutorial we cannot describe all the features of *MLwiN*, but it does have general facilities for data editing, graphing, tabulation and simple statistical summaries, all of which can be accessed through drop down menus. In addition it has a macro language, which can be used, for example, to run simulations or to carry out special purpose modelling. One of the main features is the method *MLwiN* uses to set up a model, via an ‘equation window’ in which the user specifies a model in more or less exactly the format it is usually written. Thus to specify model (1) the user would first open the equation window which, prior to any model being specified would be as follows:



(Figure 1 here)

This is the default null model with a response that is Normal with fixed predictor represented by  $X\beta$  and covariance matrix represented by  $\Omega$ . Clicking on the **N** symbol delivers a drop down menu, which allows the user to change the default distribution to binomial, Poisson or negative binomial.

Clicking on the response  $y$  allows the user to identify the response variable from a list and also the number and identification for the hierarchical levels. Clicking on the  $x_0$  term allows this to be selected from a list and also whether its coefficient  $\beta_0$  is random at particular levels of the data hierarchy. Adding a further predictor term is also a simple matter of clicking an ‘add term’ button and selecting a variable. There are simple procedures for specifying general interaction terms.

Model (1), including a random coefficient for  $X_1$  in its general form as given by (3), will be displayed in the equation window as follows:

(Figure 2 here)

Clicking on the ‘Estimates’ button will toggle the parameters between their symbolic representations and the actual estimates after a run. Likewise, the ‘Name’ button will toggle actual variable names on and off. The ‘Subs’ button allows the user to specify the form of subscripts, for example giving them names such as in the following screen where we also show the estimates and standard errors from an iterative fit:

(Figure 3 here)

In the following sections we will show some further screen shots of models and results.

### ***5. A growth data example***

We start with some simple repeated measures data and we shall use them to illustrate models of increasing complexity. The data set consists of 9 measurements made on 26 boys between the ages of 11 and 13.5 years, approximately 3 months apart [16].

Figure 4, produced by *MLwiN*, shows the mean heights by the mean age at each measurement occasion.

(Figure 4 here)

We assume that growth can be represented by a polynomial function, whose coefficients vary from individual to individual. Other functions are possible, including fractional polynomials or non-linear functions, but for simplicity we confine ourselves to examining a fourth order polynomial in age ( $t$ ) centred at an origin of 12.25 years. In some applications of growth curve modelling transformations of the time scale may be useful, often to orthogonal polynomials. In the present case the use of ordinary polynomials provides an accessible interpretation and does not lead to computational problems, for example due to near-collinearities. The model we fit can be written as follows:

$$\begin{aligned}
 y_{ij} &= \sum_{h=0}^4 \beta_{hj} t_{ij}^h + e_{ij} \\
 \beta_{0j} &= \beta_0 + u_{0j} \\
 \beta_{1j} &= \beta_1 + u_{1j} \\
 \beta_{2j} &= \beta_2 + u_{2j} \\
 \beta_{3j} &= \beta_3 \\
 \beta_{4j} &= \beta_4 \\
 \begin{pmatrix} u_0 \\ u_1 \\ u_2 \end{pmatrix} &\sim N(0, \Omega_u) \\
 \Omega_u &= \begin{pmatrix} \sigma_{u0}^2 & & \\ \sigma_{u01} & \sigma_{u1}^2 & \\ \sigma_{u02} & \sigma_{u12} & \sigma_{u2}^2 \end{pmatrix} \\
 e &\sim N(0, \sigma_e^2)
 \end{aligned} \tag{6}$$

This is a 2-level model with level 1 being ‘measurement occasion’ and level 2 ‘individual boy’. Note that we allow only the coefficients up to the second order to

vary across individuals; in the present case this provides an acceptable fit. The level 1 residual term  $e_{ij}$  represents the unexplained variation within individuals about each individual's growth trajectory. Table 1 shows the restricted maximum likelihood (REML) parameter estimates for this model. The log likelihood is calculated for the ML estimates since this is preferable for purposes of model comparison [17].

(Table 1 here)

From this table we can compute various features of growth. For example, the average growth rate (by differentiation) at age 13.25 years ( $t=1$ ) is  $6.17+2*1.13+3*0.45-4*0.38 = 8.26$  cm/year. A particular advantage of this formulation is that, for each boy, we can also estimate his random effects or 'residuals',  $u_{0j}, u_{1j}, u_{2j}$ , and use these to predict their growth curve at each age [18]. Figure 5, from *MLwiN*, shows these predicted curves (these can be produced in different colours on the screen).

(Figure 5 here)

Goldstein et al [16] show that growth over this period exhibits a seasonal pattern with growth in the Summer being about 0.5 cm greater than growth in the Winter. Since the period of the growth cycle is a year this is modelled by including a simple cosine term, which could also have a random coefficient.

In our example we have a set of individuals all of whom have 9 measurements. This restriction, however, is not necessary and (6) does not require either the same number of occasions per individual nor that measurements are made at equal intervals, since time is modelled as a continuous function. In other words we can combine data from individuals with very different measurement patterns, some of whom may only have been measured once and some who have been measured several times at irregular intervals. This flexibility, first noted by Laird and Ware [10], means that the multilevel approach to fitting repeated measures data is to be preferred to previous methods based upon a multivariate formulation assuming a common set of fixed occasions [19,20].

In these models it is assumed that the level 1 residual terms are independently distributed. We may relax this assumption, however, and in the case of repeated measures data this may be necessary, for example where measurements are taken very

close together in time. Suppose we wish to fit a model that allows for correlations between the level 1 residuals, and to start with for simplicity let us assume that these correlations are all equal. This is easily accomplished within the GLS step for the random parameters by modifying (5) to

$$\sigma_e^2 \begin{pmatrix} 1 \\ 0 \\ 1 \\ 0 \\ 0 \\ 1 \end{pmatrix} + \delta \begin{pmatrix} 0 \\ 1 \\ 0 \\ 1 \\ 1 \\ 0 \end{pmatrix} \quad (7)$$

so that the parameter  $\delta$  is the common level 1 covariance (between occasions). Goldstein et al [16] show how to model quite general nonlinear covariance functions and in particular those of the form  $\text{cov}(e_t, e_{t-s}) = \sigma_e^2 \exp(-g(\alpha, s))$ , where  $s$  is the time difference between occasions. This allows the correlation between occasions to vary smoothly as a function of their (continuous) time difference. A simple example is where  $g = \alpha s$ , which, in discrete time produces an AR(1) model. The GLS step now involves nonlinear estimation that is accomplished in a standard fashion using a Taylor series approximation within the overall iterative scheme. Pourahmadi [21,22] considers similar models but restricted to a fixed set of discrete occasions.

(Table 2 here)

Table 2 shows the results of fitting the model with  $g = \alpha s$  together with a seasonal component. If this component has amplitude, say,  $\alpha$  we can write it in the form  $\alpha \cos(t^* + \gamma)$ , where  $t^*$  is measured from the start of the calendar year. Rewriting this in the form  $\alpha_1 \cos(t^*) - \alpha_2 \sin(t^*)$  we can incorporate the  $\cos(t^*)$ ,  $\sin(t^*)$  as two further predictor variables in the fixed part of the model. In the present case  $\alpha_2$  is small and non-significant and is omitted. The results show that for measurements made three months apart the serial correlation is estimated as 0.19 ( $e^{-6.59/4}$ ) and as 0.04 ( $e^{-6.59/2}$ ) for measurements taken at six-monthly intervals. This suggests, therefore, that in practice, for such data when the intervals are no less than 6 months apart serial correlation can be ignored, but should be fitted when intervals are as small as 3 months. This will be particularly important in highly unbalanced designs where

there are some individuals with many measurements taken close together in time; ignoring serial correlation will give too much weight to the observations from such individuals.

Finally, on this topic, there will typically need to be a trade-off between modeling more random coefficients at level 2 in order to simplify or eliminate a level 1 serial correlation structure, and modeling level 2 in a parsimonious fashion so that a relatively small number of random coefficients can be used to summarise each individual. An extreme example of the latter is given by Diggle [23] who fits only a random intercept at level 2 and serial correlation at level 1.

## 6 Multivariate response data

We shall use an extension of the model for repeated measures data to illustrate how to model multivariate response data. Consider model (6) where we have data on successive occasions for each individual and in addition, for some or all individuals we have a measure, say, of their final adult height  $y_3^{(2)}$ , and their (log) income at age 25,  $y_4^{(2)}$ , where the superscript denotes a measurement made at level 2. We can include these variables as further responses by extending (6) as follows

$$\begin{aligned}
 y_{ij}^{(1)} &= \sum_{h=0}^4 \beta_{hj} t_{ij}^h + e_{ij} \\
 \beta_{0j} &= \beta_0 + u_{0j} \\
 \beta_{1j} &= \beta_1 + u_{1j} \\
 \beta_{2j} &= \beta_2 + u_{2j} \\
 \beta_{3j} &= \beta_3 \\
 \beta_{4j} &= \beta_4 \\
 y_{3j}^{(2)} &= \alpha_3 + u_{3j} \\
 y_{4j}^{(2)} &= \alpha_4 + u_{4j} \\
 \begin{pmatrix} u_0 \\ u_1 \\ u_2 \\ u_3 \\ u_4 \end{pmatrix} &\sim N(0, \Omega_u) \\
 e &\sim N(0, \sigma_e^2)
 \end{aligned} \tag{8}$$

We now have a model where there are response variables defined at level 1 (with superscript (1)) and also at level 2 (with superscript (2)). For the level 2 variables we have specified only an intercept term in the fixed part, but quite general functions of individual level predictors, such as gender, are possible. The level 2 responses have no component of random variation at level 1 and their level 2 residuals covary with the polynomial random coefficients from the level 1 repeated measures response.

The results of fitting this model allow us to quantify the relationships between growth events, such as growth acceleration (differentiating twice) at  $t=0$ , age 12.25 years, ( $2\beta_{2j}$ ) and adult height and also to use measurements taken during the growth period to make efficient predictions of adult height or income. We note that for individual  $j$

the estimated (posterior) residuals  $\hat{u}_{3j}, \hat{u}_{4j}$  are the best linear unbiased predictors of the individual's adult values; where we have only a set of growth period measurements for an individual these therefore provide the required estimates. Given the set of model parameters, therefore, we immediately obtain a system for efficient adult measurement prediction given a set of growth measurements [24].

Suppose, now, that we have no growth period measurements and just the two adult measurements for each individual. Model (8) reduces to

$$\begin{aligned}
 y_{3j}^{(2)} &= \alpha_3 + u_{3j} \\
 y_{4j}^{(2)} &= \alpha_4 + u_{4j} \\
 \begin{pmatrix} u_3 \\ u_4 \end{pmatrix} &\sim N(0, \Omega_u) \\
 \mathbf{V}_{1j} &= 0, \mathbf{V}_{2j} = \begin{pmatrix} \sigma_{u3}^2 & \\ \sigma_{u34} & \sigma_{u4}^2 \end{pmatrix}
 \end{aligned} \tag{9}$$

So that we can think of this as a 2 level model with no level 1 variation and every level 2 unit containing just two level 1 units. The explanatory variables for the simple model given by (9) are just two dummy variables defining, alternately the two responses. Thus we can write (9) in the more compact general form

$$\begin{aligned}
 y_{ij} &= \sum_{h=1}^2 \beta_{0hj} x_{hij}, \quad x_{1ij} = \begin{cases} 1 & \text{if response 1} \\ 0 & \text{if response 2} \end{cases}, \quad x_{2ij} = 1 - x_{1ij} \\
 \beta_{0hj} &= \beta_{0h} + u_{hij} \\
 \begin{pmatrix} u_1 \\ u_2 \end{pmatrix} &= \begin{pmatrix} \sigma_{u1}^2 & \\ \sigma_{u12} & \sigma_{u2}^2 \end{pmatrix}
 \end{aligned} \tag{10}$$

Note that there is no need for every individual to have both responses and so long as we can consider 'missing' responses as random, the IGLS algorithm will supply maximum likelihood estimates. We can add further covariates to the model in a straightforward manner by forming interactions between them and the dummy variables defining the separate response intercepts.

The ability to fit a multivariate linear model with randomly missing responses finds a number of applications, for example where matrix or rotation designs are involved [18, Chapter 4], each unit being allocated, at random, a subset of responses. The

possibility of having additionally level 1 responses allows this to be used as a very general model for meta analysis where there are several studies (level 2 units) for some of which responses are available only in summary form at level 2 and for others detailed level 1 responses are available. Goldstein et al. [25] provide a detailed example.

## 6. Cross classified and multiple membership structures

Across a wide range of disciplines it is commonly the case that data have a structure that is not purely hierarchical. Individuals may be clustered not only into hierarchically ordered units (for example occasions nested within patients nested within clinics), but may also belong to more than one type of unit at a given level of a hierarchy. Consider the example of a livestock animal such as a cow where there are a large number of mothers, each producing several female offspring that are eventually reared for milk on different farms. Thus, an offspring might be classified as belonging to a particular combination of mother and farm, in which case they will be identified by a *cross classification* of these.

Raudenbush [26] and Rasbash and Goldstein [27] present the general structure of a model for handling complex hierarchical structuring with random cross classifications. For example, assuming that we wish to formulate a linear model for the milk yield of offspring taking into account both the mother and the farm, then we have a cross classified structure, which can be modelled as follows:

$$y_{i(j_1 j_2)} = (X\beta)_{i(j_1 j_2)} + u_{j_1} + u_{j_2} + e_{i(j_1 j_2)}, \quad (11)$$

$$j_1 = 1, \dots, J_1, \quad j_2 = 1, \dots, J_2, \quad i = 1, \dots, N$$

in which the yield of offspring  $i$ , belonging to the combination of mother  $j_1$  and farm  $j_2$ , is predicted by a set of fixed coefficients  $(X\beta)_{i(j_1, j_2)}$ . The random part of the model is given by two level 2 residual terms, one for the mother ( $u_{j_1}$ ) and one for the farm ( $u_{j_2}$ ), together with the usual level 1 residual term for each offspring. Decomposing the variation in such a fashion allows us to see how much of it is due to the different classifications. This particular example is somewhat oversimplified, since we have ignored paternity and we would also wish to include factors such as age of mother, parity of offspring etc. An application of this kind of modelling to a more



complex structure involving Salmonella infection in chickens is given by Rasbash and Browne [28].

Considering now just the farms, and ignoring the mothers, suppose that the offspring often change farms, some not at all and some several times. Suppose also that we know, for each offspring, the weight  $w_{ij_2}$ , associated with the  $j_2$ -th farm for offspring

$i$  with  $\sum_{j_2=1}^{J_2} w_{ij_2} = 1$ . These weights, for example, may be proportional to the length of

time an offspring stays in a particular farm during the course of our study. Note that we allow the possibility that for some (perhaps most) animals only one farm is involved so that one of these probabilities is one and the remainder are zero. Note that when all level 1 units have a single non-zero weight of 1 we obtain the usual purely hierarchical model. We can write for the special case of membership of up to two farms  $\{1,2\}$ :

$$\begin{aligned} y_{i(1,2)} &= (X\beta)_{i(1,2)} + w_{i1}u_1 + w_{i2}u_2 + e_{i(1,2)} \\ w_{i1} + w_{i2} &= 1 \end{aligned} \quad (12)$$

and more generally:

$$\begin{aligned} y_{i\{j\}} &= (X\beta)_{i\{j\}} + \sum_{h \in \{j\}} w_{ih}u_h + e_{i\{j\}} \\ \sum_h w_{ih} &= 1, \quad \text{var}(u_h) = \sigma_u^2 \\ \text{var}\left(\sum_h w_{ih}u_h\right) &= \sigma_u^2 \sum_h w_{ih}^2, \end{aligned} \quad (13)$$

Thus, in the particular case of membership of just two farms with equal weights we have

$$w_{i1} = w_{i2} = 0.5, \quad \text{var}\left(\sum_h w_{ih}u_h\right) = \sigma_u^2 / 2$$

Further details of this model are given by Hill and Goldstein [29].

An extension of the multiple membership model is also possible and has important applications, for example in modelling spatial data. In this case we can write

$$\begin{aligned}
y_{i\{j_1\}\{j_2\}} &= (X\beta)_{i\{j\}} + \sum_{h \in \{j_1\}} w_{1ih} u_{1h} + \sum_{h \in \{j_2\}} w_{2ih} u_{2h} + e_{i\{j\}} \\
\sum_h w_{1ih} &= W_1, \quad \sum_h w_{2ih} = W_2, \quad \text{var}(u_{1h}) = \sigma_{u1}^2 \quad \text{var}(u_{2h}) = \sigma_{u2}^2 \\
\text{cov}(u_{1h}, u_{2h}) &= \sigma_{u12}, \quad j = \{j_1, j_2\}
\end{aligned} \tag{14}$$

There are now two sets of higher level units that influence the response and in general we can have more than two such sets. In spatial models one of these sets is commonly taken to be the area where an individual (level 1) unit occurs and so does not have a multiple membership structure (since each individual belongs to just one area, that is we replace  $\sum_h w_{1ih} u_{1h}$  by  $u_{1j_1}$ ). The other set consists of those neighbouring units that are assumed to have an effect. The weights will need to be carefully chosen; in spatial models  $W_2$  is typically chosen to be equal 1 (see Langford et al. [30] for an example). Another application for a model such as (14) is for household data where households share facilities, for example an address. In this case the household that an individual resides in will belong to one set and the other households at the address will belong to the other set. Goldstein et al. [31] give an application of this model to complex household data.

## 7. Meta analysis

Meta analysis involves the pooling of information across studies in order to provide both greater efficiency for estimating treatment effects and also for investigating why treatments effects may vary. By formulating a general multilevel model we can do both of these efficiently within a single model framework, as has already been indicated and was suggested by several authors [32, 33]. In addition we can combine data that are provided at either individual subject level or aggregate level or both. We shall look at a simple case but this generalises readily [25].

Consider an underlying model for individual level data where a pair of treatments are being compared and results from a number of studies or centres are available. We write a basic model, with a continuous response  $Y$  as

$$\begin{aligned}
y_{ij} &= (X\beta)_{ij} + \beta_2 t_{ij} + u_j + e_{ij} \\
\text{var}(u_j) &= \sigma_u^2, \quad \text{var}(e_{ij}) = \sigma_e^2
\end{aligned} \tag{15}$$

with the usual assumptions of Normality etc. The covariate function is designed to adjust for initial clinic and subject conditions. The term  $t_{ij}$  is a dummy variable defining the treatment (0 for treatment A, 1 for treatment B). The random effect  $u_j$  is a study effect and the  $e_{ij}$  are individual level residuals. Clearly this model can be elaborated in a number of ways, by including random coefficients at level 2 so that the effect of treatment varies across studies, and by allowing the level 1 variance to depend on other factors such as gender or age.

Suppose now that we do not have individual data available but only means at the study level. If we average (15) to the study level we obtain

$$y_j = (X\beta)_{.j} + \beta_2 t_{.j} + u_j + e_j \quad (16)$$

where  $y_j$  is the mean response for the  $j$ -th study etc. The total residual variance for study  $j$  in this model is  $\sigma_u^2 + \sigma_e^2 / n_j$  where  $n_j$  is the size of the  $j$ -th study. It is worth noting at this point that we are ignoring, for simplicity, levels of variation that might exist within studies, such as that between sites for a multi-site study. If we have the values of  $y_j$ ,  $(X\beta)_{.j}$ ,  $t_{.j}$  where the latter is simply the proportion of subjects with treatment B in the  $j$ -th study, and also the value of  $n_j$  then we will be able to obtain estimates for the model parameters, so long as the  $n_j$  differ. Such estimates, however, may not be very precise and extra information, especially about the value of  $\sigma_e^2$  will improve them.

Model (16) therefore forms the basis for the multilevel modelling of aggregate level data. In practice the results of studies will often be reported in non-standard form, for example with no estimate of  $\sigma_e^2$  but it may be possible to estimate this from reported test statistics. In some cases, however, the reporting may be such that the study cannot be incorporated in a model such as (16). Goldstein et al [25] set out a set of minimum reporting conventions for meta analysis studies subsequently to be carried out.

While it is possible to perform a meta analysis with only aggregate level data, it is clearly more efficient to utilise individual level data where these are available. In general, therefore, we will need to consider models that have mixtures of individual

and aggregate data, even perhaps within the same study. We can do this straightforwardly by specifying a model which is just the combination of (15) and (16), namely

$$\begin{aligned}
y_{ij} &= \beta_0 + \beta_1 x_{ij} + \beta_2 t_{ij} + u_j + e_{ij} \\
y_j &= \beta_0 + \beta_1 x_j + \beta_2 t_j + u_j + e_j z_j \\
z_j &= \sqrt{n_j^{-1}}, \quad e_j \equiv e_{ij}
\end{aligned} \tag{17}$$

What we see is that the common level 1 and level 2 random terms link together the separate models and allow a joint analysis that makes fully efficient use of the data. Several issues immediately arise from (17). One is that the same covariates are involved. This is also a requirement for the separate models. If some covariate values are missing at either level then it is possible to use an imputation technique to obtain estimates, assuming a suitable random missingness mechanism. The paper by Goldstein et al [25] discusses generalisations of (17) for several treatments and the procedure can be extended to generalised linear models.

## 8. Generalised linear models

So far we have dealt with linear models, but all of those so far discussed can be modified using nonlinear link functions to give generalized linear multilevel models. We shall not discuss these in detail (see Goldstein [18] for details and some applications) but for illustration we will describe a 2 level model with a binary response.

Suppose the outcome of patients in intensive care is recorded simply as survived (0) or died (1) within 24 hours of admission. Given a sample of patients from a sample of intensive care units we can write one model for the probability of survival as

$$\begin{aligned}
\text{logit}(\pi_{ij}) &= (X\beta)_{ij} + u_j \\
y_{ij} &\sim \text{Bin}(\pi_{ij}, 1)
\end{aligned} \tag{18}$$

Equation (18) uses a standard logit link function assuming binomial (Bernoulli) error distribution for the (0,1) response  $y_{ij}$ . The level 2 random variation is described by the term  $u_j$  within the linear predictor. The general interpretation is similar to that for a continuous response model, except that the level 1 variation is now a function of the

predicted value  $\pi_{ij}$ . While in (18) there is no separate estimate for the level 1 variance, we may wish to fit extra-binomial variation which will involve a further parameter.

We can modify (18) using alternate link functions, for example the logarithm if the response is a count and can allow further random coefficients at level 2. The response can be a multcategory variable, either ordered or unordered and this provides an analogue to the multivariate models for continuously distributed responses. As an example of an ordered categorical response consider extending the previous outcome to 3 categories; survival without impairment (1), survival with impairment (2), death (3). If  $\pi_{ij}^{(h)}$ ,  $h = 1, 2, 3$  are respectively the probabilities for each of these categories we can write a *proportional odds* model using a logit link as

$$\begin{aligned} \text{logit}(\pi_{ij}^{(1)}) &= \alpha^{(1)} + (X\beta)_{ij} + u_j^{(1)} \\ \text{logit}(\pi_{ij}^{(1)} + \pi_{ij}^{(2)}) &= \alpha^{(2)} + (X\beta)_{ij} + u_j^{(2)} \end{aligned} \tag{19}$$

and the set of three (0,1) observed responses for each patient is assumed to have a multinomial distribution with mean vector given by  $\pi_{ij}^{(h)}$ ,  $h = 1, 2, 3$ . Since the probabilities add to 1 we require two lines in (19) which differ only in terms of the overall level given by the intercept term as well as allowing for these to vary across units.

Unlike in the continuous Normal response case, maximum likelihood estimation is not straightforward and beyond fairly simple 2 level models involves a considerable computational load typically using numerical integration procedures [34]. For this reason approximate methods have been developed based upon series expansions and using quasilielihood approaches [18] which perform well under a wide range of circumstances but can break down in certain conditions, especially when data are sparse with binary responses. High order Laplace approximations have been found to perform well [35] as have simulation-based procedures such as MCMC (see below).

It is worth mentioning one particular situation where care is required in using generalized linear multilevel models. In (15) and (16) we assume that the level 1 responses are independently distributed with finite probabilities distributed according to the specified model. In some circumstances such an assumption may not be sensible. Consider a repeated measures model where the health status (satisfactory/not satisfactory) of individuals is repeatedly assessed. Some individuals will typically always respond ‘satisfactory’ whereas some others can be expected to respond always ‘not satisfactory’. For these individuals the underlying probabilities are either zero or one, which violates the model assumptions and what one finds if one tries to fit a model where there are non-negligible numbers of such individuals are noticeable amounts of underdispersion. Barbosa and Goldstein [36] discuss this problem and propose a solution based upon fitting a serial correlation structure.

We can also have multivariate multilevel models with mixtures of discrete and continuous responses. Certain of these can be fitted in MlwiN using quasilielihood procedures [12] and MCMC procedures for such models are currently being implemented. See also Olsen and Shaffer [37] for an alternative approach.

## 9. Survival models

Several different formulations for survival data modeling are available; to illustrate how these can be extended to multilevel structures, where they are often referred to as frailty models [38], we consider the proportional hazards (Cox) model and a piecewise, discrete time, model. Goldstein [18] gives other examples.

Consider a simple 2-level model with, say, patients within hospitals or occasions within subjects. As in the standard single level case we consider each time point in the data as defining a *block* indicated by  $l$  at which some observations come to the end of their duration due to either failure or censoring and some remain to the next time or block. At each block there is therefore a set of observations - the total *risk set*. To illustrate how the model is set up we can think of the data sorted so that each observation within a block is a level-1 unit, above which, in the repeated measures case, there are occasions at level-2 and subjects at level 3. The ratio of the hazard for

the unit which experiences a failure at a given occasion referred to by  $(j, k)$  to the sum of the hazards of the remaining risk set units [39] is

$$\frac{\exp(\beta_1 x_{1ijk} + u_{jk})}{\sum_{j,k} \exp(\beta_1 x_{1ijk} + u_{jk})} \quad (20)$$

where  $j$  and  $k$  refer to the *real* levels 2 and 3, for example occasion and subject. At each block denoted by  $l$  the response variable may be defined *for each member of the risk set* as

$$y_{ijk(l)} = \begin{cases} 1 & \text{failed} \\ 0 & \text{not} \end{cases}$$

Because of equivalence between the likelihood for the multinomial and Poisson distributions, the latter is used to fit model (20). This can be written as

$$E(y_{ijk(l)}) = \exp(\alpha_l + X_{jk} \beta_k) \quad (21)$$

Where there are ties within a block then more than one response will be non-zero. The terms  $\alpha_l$  fit the underlying hazard function as a '*blocking factor*', and can be estimated by fitting either a set of parameters, one for each block, or a smoothed polynomial curve over the blocks numbered  $1, \dots, p$ . Thus if the  $h^{th}$  block is denoted by  $h$ ,  $\alpha_l$  is replaced by a low order polynomial, order  $m$ ,  $\sum_{t=0}^m \gamma_t h^t$ , where the  $\gamma_t$  are (nuisance) parameters to be estimated.

Having set up this model, the data are now sorted into the real 2-level structure, for example in the repeated measures case by failure times within subjects with occasions within the failure times. This retains proportional hazards within subjects. In this formulation the Poisson variation is defined at level-1, there is no variation at level-2 and the between-subject variation is at level 3. Alternatively we may wish to preserve *overall* proportionality, in which case the failure times define level 3 with no variation at that level. See Goldstein [18] for a further discussion of this.

Consider now a piecewise survival model. Here the total time interval is divided into short intervals during which the probability of failure, given survival up to that point, is assumed constant. Denote these intervals by  $t$  ( $1, 2, \dots, T$ ) so that the hazard at time  $t$  is the probability that, given survival up to the end of time interval  $t-1$ , failure occurs in the next interval. At the start of each interval we have a ‘risk set’  $n_t$  consisting of the survivors and during the interval  $r_t$  fail. If censoring occurs during interval  $t$  then this observation is removed from that interval (and subsequent ones) and does not form part of the risk set. A simple, *single level*, model for the probability can be written as

$$\pi_{i(t)} = f[\alpha_t z_{it}, (\beta X)_{it}] \quad (22)$$

where  $z_t = \{z_{it}\}$  is a dummy variable for the  $t$ -th interval and  $\alpha_t$ , as before, is a ‘blocking factor’ defining the underlying hazard function at time  $t$ . The second term is a function of covariates. A common formulation would be the logit model and a simple such model in which the first blocking factor has been absorbed into the intercept term could be written as

$$\text{logit}(\pi_{i(t)}) = \beta_0 + \alpha_t z_{it} + \beta_1 x_{1i}, \quad (z_2, z_3, \dots, z_T) \quad (23)$$

Since the covariate varies across individuals, in general the data matrix will consist of one record for each individual within each interval, with a (0,1) response indicating survival or failure. The model can now be fitted using standard procedures, assuming a binomial error distribution. As before, instead of fitting  $T-1$  blocking factors, we can fit a low order polynomial to the sequentially numbered time indicator. The logit function can be replaced by, for example, the complementary log-log function that gives a proportional hazards model, or, say, the probit function and note that we can incorporate time-varying covariates such as age.

For the extension to a two level model we write

$$\text{logit}(\pi_{ij(t)}) = \beta_0 + \sum_{h=1}^p \alpha_h^* (z_{it}^*)^h + \beta_1 x_{1ij} + u_j \quad (24)$$

where  $u_j$  is the ‘effect’, for example, of the  $j$ -th clinic, and is typically assumed to be distributed Normally with zero mean and variance  $\sigma_u^2$ . We can elaborate (24) using random coefficients, resulting in a heterogeneous variance structure, further levels of nesting etc. This is just a 2-level binary response model and can be fitted as described



earlier. The data structure has two levels so that individuals will be grouped (sorted) within clinics. For a competing risks model, with more than one outcome we can use the 2-level formulation for a multcategory response described above. The model can be used with repeated measures data where there are repeated survival times within individuals, for example multiple pregnancy states.

## **Bayesian modelling**

So far we have considered the classical approach to fitting multilevel models. If we add prior distributional assumptions to the parameters of the models so far considered we can fit the same range of models from a Bayesian perspective, and in most applications this will be based upon MCMC methods. A detailed comparison of Bayesian and likelihood procedures for fitting multilevel models is given in Browne and Draper [40]. A particular advantage of MCMC methods is that they yield inferences based upon samples from the full posterior distribution and allow exact inference in cases where, as mentioned above, the likelihood based methods yield approximations.

Due also to their approach of generating a sample of points from the full posterior distributions they can give accurate interval estimates for non Gaussian parameter distributions.

In MCMC sampling we are interested in generating samples of values from the joint posterior distribution of all the unknown parameters rather than finding the maximum of this distribution. Generally it is not possible to generate directly from this joint distribution, so instead the parameters are split into groups and for each group in turn we generate a set of values from its conditional posterior distribution. This can be shown to be equivalent to sampling directly from the joint posterior distribution.

There are two main MCMC procedures that are used in practice, Gibbs sampling [41] and Metropolis-Hastings (MH) [42,43] sampling. When the conditional posterior for a group of parameters has a standard form for example a Normal distribution then we can generate values from it directly and this is known as Gibbs sampling. When the distribution is not of standard form then it may still be possible to use Gibbs sampling by constructing the distribution using forms of Gibbs sampling such as Adaptive rejection sampling [44].

The alternative approach is to use MH sampling where values are generated from another distribution called a proposal distribution rather than the conditional posterior distribution. These values are then either accepted or rejected in favour of the current values by comparing the posterior probabilities of the joint posterior at the current and proposed new values. The acceptance rule is designed so that MH is effectively sampling from the conditional posterior even though we have used an arbitrary proposal distribution; nevertheless, choice of this proposal distribution is important for efficiency of the algorithm.

MCMC algorithms produce chains of serially correlated parameter estimates and consequently often have to be run for many iterations to get accurate estimates. Many diagnostics are available to gauge approximately for how long to run the MCMC methods. The chains are also started from arbitrary parameter values and so it is common practise to ignore the first  $N$  iterations (known as a burn-in period) to allow the chains to move away from the starting value and settle at the parameters equilibrium distribution.

We give here an outline of the Gibbs sampling procedure for fitting a general Normal two level model.

$$\begin{aligned}
 y_{ij} &= \mathbf{X}_{ij}\boldsymbol{\beta} + \mathbf{Z}_{ij}u_j + e_{ij} \\
 u_j &\sim N(0, \boldsymbol{\Omega}_u), e_{ij} \sim N(0, \sigma_e^2), \\
 i &= 1, \dots, n_j, j = 1, \dots, J, \sum_j n_j = N
 \end{aligned}$$

We will include generic conjugate prior distributions for the fixed effects and variance parameters as follows:

$$\boldsymbol{\beta} \sim N(\boldsymbol{\mu}_p, S_p), \boldsymbol{\Omega}_u \sim W^{-1}(\boldsymbol{v}_u, S_u), \sigma_e^2 \sim SI\chi^2(\boldsymbol{v}_e, s_e^2)$$

The Gibbs sampling algorithm then involves simulating from the following four sets of conditional distributions

$$p(\boldsymbol{\beta} | y, \sigma_e^2, u) \sim N\left(\frac{\hat{D}}{\sigma_e^2} \left( \sum_{j=1}^J \sum_{i=1}^{n_j} \mathbf{X}_{ij}^T (y_{ij} - \mathbf{Z}_{ij} u_j) + S_p^{-1} \boldsymbol{\mu}_p \right), \hat{D}\right)$$

$$p(u_j | y, \Omega_u, \sigma_e^2, \boldsymbol{\beta}) \sim N\left(\frac{\hat{D}_j}{\sigma_e^2} \sum_{i=1}^{n_j} \mathbf{Z}_{ij}^T (y_{ij} - \mathbf{X}_{ij} \boldsymbol{\beta}), \hat{D}_j\right)$$

$$p(\Omega_u | u) \sim W^{-1}\left(J + v_u, \sum_{j=1}^J u_j u_j^T + S_u\right)$$

$$p(\sigma_e^2 | y, \boldsymbol{\beta}, u) \sim \Gamma^{-1}\left(\frac{N + v_e}{2}, \frac{1}{2}(v_e s_e^2 + \sum_{j=1}^J \sum_{i=1}^{n_j} e_{ij}^2)\right)$$

$$\text{where } \hat{D} = \left( \sum_{ij} \frac{\mathbf{X}_{ij}^T \mathbf{X}_{ij}}{\sigma_e^2} + S_p^{-1} \right)^{-1} \text{ and } \hat{D}_j = \left( \frac{1}{\sigma_e^2} \sum_{i=1}^{n_j} \mathbf{Z}_{ij}^T \mathbf{Z}_{ij} + \Omega_u^{-1} \right)^{-1}.$$

Note that in this algorithm we have used generic prior distributions. This allows the incorporation of informative prior information but generally we will not have this information and so will use so-called ‘diffuse’ prior distributions that reflect our lack of knowledge. Since there are only 26 level 2 units from which we are estimating a 3 x 3 covariance matrix, the exact choice of prior is important. We here use the following set of priors for the child growth example considered in section 4.

$$p(\boldsymbol{\beta}_0) \propto 1, \quad p(\boldsymbol{\beta}_1) \propto 1, \quad p(\boldsymbol{\beta}_2) \propto 1, \quad p(\boldsymbol{\beta}_3) \propto 1, \quad p(\boldsymbol{\beta}_4) \propto 1, \quad p(\boldsymbol{\beta}_5) \propto 1$$

$$p(\Omega_u) \sim \text{inverse Wishart}_3[3, 3 * S_u], \quad S_u = \begin{bmatrix} 64.0 & & \\ 8.32 & 2.86 & \\ 1.42 & 0.92 & 0.67 \end{bmatrix}$$

$$p(\sigma_{e0}^2) \sim \text{inverse gamma}(0.001, 0.001)$$

The inverse Wishart prior matrix is based upon the REML estimates, chosen to be ‘minimally informative’, with degrees of freedom equal to the order of the covariance matrix.

The results of running model 5 for 50,000 iterations after a burn-in of 500 can be seen in table 3. Here we see that the fixed effect estimates are very similar to the estimates obtained by the maximum likelihood method. The variance (chain mean) estimates are however inflated due to the skewness of the variance parameters. Modal estimates of the variance parameters, apart from that for the quadratic coefficient, are closer as is to be expected. If we had used the uniform prior  $p(\Omega_u) \propto 1$  for the covariance

matrix the estimates of the fixed coefficients are little changed but the covariance matrix estimates are noticeably different. For example the variance associated with the intercept is now 93.0 and those for the linear and quadratic coefficients become 4.2 and 1.1 respectively.

(Table 3 here)

Figure 6 shows the diagnostic screen produced by *MLwiN* following this MCMC run.

(Figure 6 here)

The top left hand box shows the trace for this parameter, the level two intercept variance. This looks satisfactory and this is confirmed by the estimated autocorrelation function (ACF) and partial autocorrelation function (PACF) below. A kernel density estimate is given at the top right and the bottom left box is a plot of the Monte Carlo standard error against number of iterations in the chain. The summary statistics give quantiles, mean and mode together with accuracy diagnostics that indicate the required chain length.

The MCMC methods are particularly useful in models like the cross classified and multiple membership models discussed in section 6. This is because, whereas the maximum likelihood methods involve constructing large constrained variance matrices for these models, the MCMC methods simulate conditional distributions in turn and so do not have to adjust to the structure of the model.

For model fitting, one strategy (but not the only one) is to use the maximum or quasi-likelihood methods for performing model exploration and selection due to speed. Then MCMC methods could be used for inference on the final model to obtain accurate interval estimates.

## **Bootstrapping**

Like MCMC the bootstrap allows inferences to be based upon (independent) chains of values and can be used to provide exact inferences and corrections for bias. Two forms of bootstrapping have been studied to date, parametric bootstrapping, especially

for correcting biases in generalised linear models [45], and nonparametric bootstrapping based upon estimated residuals [46]. The fully parametric bootstrap for a multilevel model works as for a single level model with simulated values generated from the estimated covariance matrices at each level. Thus, for example, for model (2) each bootstrap sample is created using a set of  $u_j, e_{ij}$  sampled from  $N(0, \sigma_u^2), N(0, \sigma_e^2)$  respectively.

In the nonparametric case, full ‘unit resampling’ is generally only possible by resampling units at the highest level. For generalised linear models, however, we can resample posterior residuals, once they have been adjusted to have the correct (estimated) covariance structures and this can be shown to possess particular advantages over a fully parametric bootstrap where asymmetric distributions are involved [46].

## **In Conclusion**

The models that we have described in this paper represent a powerful set of tools available to the data analyst for exploring complex data structures. They are being used in many areas, including health, with great success in providing insights that are unavailable with more conventional methods. There is a growing literature extending these models, for example to multilevel structural equation models, and especially to the application of the multiple membership models in areas such as population demography [31]. An active email discussion group exists which welcomes new members ([www.jiscmail.ac.uk/multilevel](http://www.jiscmail.ac.uk/multilevel)). The data set used for illustration in this paper is available on the Centre for Multilevel Modelling website ([multilevel.ioe.ac.uk/download/macros.html](http://multilevel.ioe.ac.uk/download/macros.html)) as an *MLwiN* worksheet.

## References

1. Snijders, T. and Bosker, R. (1999). *Multilevel Analysis*. London, Sage:
2. Leyland, A. H. and Goldstein, H. (Editors). (2001). *Multilevel Modelling of Health Statistics*. Chichester, Wiley.
3. Sullivan, L. M., Dukes, K. A. and Losina, E. (1999). Tutorial in biostatistics: an introduction to hierarchical linear modelling. *Statistics in Medicine* **18**: 855-888.
4. Zeger, S. L., Liang, K. and Albert, P. (1988). Models for longitudinal data: A generalised estimating equation approach. *Biometrics* **44**: 1049-1060.
5. Liang, K.-Y., Zeger, S. L. and Qaqish, B. (1992). Multivariate regression analyses for categorical data. *J.Royal Statist.Soc.B.* **54**: 3-40.
6. Heagerty, P. J. and Zeger, S. L. (2000). Marginalized multilevel models and likelihood inference (with discussion). *Statistical Science* **15**: 1-26.
7. Lindsey, J. K. and Lambert, P. (1998). On the appropriateness of marginal models for repeated measurements in clinical trials. *Statistics in Medicine* **17**: 447-469.
8. Goldstein, H. and Rasbash, J. (1992). Efficient computational procedures for the estimation of parameters in multilevel models based on iterative generalised least squares. *Computational Statistics and Data Analysis* **13**: 63-71.
9. Longford, N. T. (1987). A fast scoring algorithm for maximum likelihood estimation in unbalanced mixed models with nested effects. *Biometrika* **74**: 812-27.
10. Laird, N. M. and Ware, J. H. (1982). Random-effects models for longitudinal data. *Biometrics* **38**: 963-974.
11. Bryk, A. S. and Raudenbush, S. W. (1992). *Hierarchical Linear Models*. Newbury Park, California, Sage:
12. Rasbash, J., Browne, W., Goldstein, H., Yang, M., et al. (2000). *A user's guide to MlwiN (Second Edition)*. London, Institute of Education:
13. Goldstein, H. (1986). Multilevel mixed linear model analysis using iterative generalised least squares. *Biometrika* **73**: 43-56.
14. Zhou, X., Perkins, A. J. and Hui, S. L. (1999). Comparisons of software packages for generalized linear multilevel models. *American Statistician* **53**: 282-290.

15. Spiegelhalter, D. J., Thomas, A. and Best, N. G. (2000). *WinBUGS Version 1.3: User Manual*. Cambridge, MRC Biostatistics Research Unit:
16. Goldstein, H., Healy, M. J. R. and Rasbash, J. (1994). Multilevel time series models with applications to repeated measures data. *Statistics in Medicine* **13**: 1643-55.
17. Kenward, M. G. and Roger, J. H. (1997). Small sample inference for fixed effects from restricted maximum likelihood. *Biometrics* **53**: 983-997.
18. Goldstein, H (1995). *Multilevel Statistical Models*. London, Arnold
19. Grizzle, J. C. and Allen, D. M. (1969). An analysis of growth and dose response curves. *Biometrics* **25**: 357-61.
20. Albert, P. S. (1999). Longitudinal data analysis (repeated measures) in clinical trials. *Statistics in Medicine* **18**: 1707-1732.
21. Pourahmadi, M. (1999). Joint mean-covariance models with applications to longitudinal data: unconstrained parameterisation. *Biometrika* **86**: 677-690.
22. Pourahmadi, M. (2000). Maximum likelihood estimation of generalised linear models for multivariate Normal covariance matrix. *Biometrika* **87**: 425-436.
23. Diggle, P. J. (1988). An approach to the analysis of repeated measurements. *Biometrics* **44**: 959-971.
24. Goldstein, H. (1989). Flexible Models for the analysis of growth data with an application to height prediction. *Rev.Epidem.et Sante Public* **37**: 477-484.
25. Goldstein, H., Yang, M., Omar, R., Turner, R., et al. (2000). Meta analysis using multilevel models with an application to the study of class size effects. *Journal of the Royal Statistical Society, Series C* **49**: 399-412.
26. Raudenbush, S. W. (1993). A crossed random effects model for unbalanced data with applications in cross sectional and longitudinal research. *Journal of Educational Statistics* **18**: 321-349.
27. Rasbash, J. and Goldstein, H. (1994). Efficient analysis of mixed hierarchical and cross classified random structures using a multilevel model. *Journal of Educational and Behavioural statistics* **19**: 337-50.

28. Rasbash, J. and Browne, W. (2001). Non-hierarchical multilevel models. In *Multilevel Modelling of Health Statistics*. A. Leyland and H. Goldstein (Eds.). Chichester, Wiley.
29. Hill, P. W. and Goldstein, H. (1998). Multilevel modelling of educational data with cross-classification and missing identification of units. *Journal of Educational and Behavioural statistics* **23**: 117-128.
30. Langford, I., Leyland, A. H., Rasbash, J. and Goldstein, H. (1999). Multilevel modelling of the geographical distributions of diseases. *Journal of the Royal Statistical Society, Series C* **48**: 253-68.
31. Goldstein, H., Rasbash, J., Browne, W., Woodhouse, G., et al. (2000). Multilevel models in the study of dynamic household structures. *European Journal of Population* **16**, 373-38732.
32. Hedges, L. V. and Olkin, I. O. (1985). *Statistical methods for meta analysis*. Orlando, Florida, Academic Press.
33. Raudenbush, S. and Bryk, A. S. (1985). Empirical Bayes meta-analysis. *Journal of Educational Statistics* **10**: 75-98.
34. Hedeker, D. and Gibbons, R. D. (1994). A random effects ordinal regression model for multilevel analysis. *Biometrics* **50**: 933-44.
35. Raudenbush, S. W., Yang, M. and Yosef, M. (2000). Maximum likelihood for generalised linear models with nested random effects via high-order multivariate Laplace approximation. *Journal of Computational and Graphical Statistics*. **9**: 141-157.
36. Barbosa, M. F. and Goldstein, H. (2000). Discrete response multilevel models for repeated measures; an application to voting intentions data. *Quality and Quantity* **34**: 323-330.
37. Olsen, M. K. and Schafer, J. L. (2001). A two-part random effects model for semi-continuous longitudinal data. *Journal of the American Statistical Association* **96**: 730-745.
38. Clayton, D. G. (1991). A monte carlo method for Bayesian inference in frailty models. *Biometrics* **47**: 467-85.



39. McCullagh, P. and Nelder, J. (1989). *Generalised linear models*. London, Chapman and Hall:
40. Browne, W. and Draper, D. (2000). Implementation and performance issues in the Bayesian and likelihood fitting of multilevel models. *Computational Statistics*, *15*, 391-420
41. Geman, S. and Geman, D. (1984). Stochastic relaxation, Gibbs distributions and the Bayesian restoration of images. *IEEE Transactions on Pattern Analysis and Machine Intelligence* **45**: 721-741.
42. Metropolis, N., Rosenbluth, A.W., Rosenbluth, M.N., Teller, A.H., and Teller, E. (1953). Equations of state calculations by fast computing machines. *Journal of Chemical Physics* **21**: 1087-1092.
43. Hastings, W.K. (1970). Monte Carlo sampling methods using Markov chains and their applications. *Biometrika* **57**: 97-109.
44. Gilks, W.R. and Wild, P. (1992). Adaptive rejection sampling for Gibbs sampling. *Journal of the Royal Statistics Society, Series C* **41**, 337-348.
45. Goldstein, H. (1996). Consistent estimators for multilevel generalised linear models using an iterated bootstrap. *Multilevel Modelling Newsletter* **8**(1): 3-6.
46. Carpenter, J., Goldstein, H. and Rasbash, J. (1999). A nonparametric bootstrap for multilevel models. *Multilevel modelling Newsletter*, **11** (1), 2-5.

## Figures

Figure 1. Default equation screen with model unspecified.

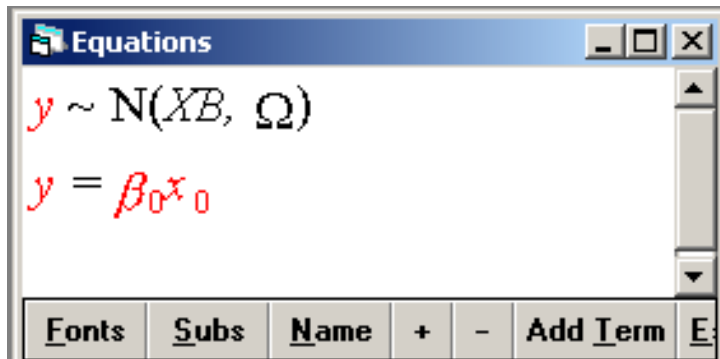


Figure 2. Equation screen with model display.

The screenshot shows a software window titled "Equations" with a scrollable text area containing the following mathematical expressions:

$$y_{ij} \sim N(XB, \Omega)$$
$$y_{ij} = \beta_{0ij}x_{0i} + \beta_{1ij}x_{1ij}$$
$$\beta_{0ij} = \beta_0 + u_{0ij} + e_{0ij}$$
$$\beta_{1ij} = \beta_1 + u_{1ij}$$
$$\begin{bmatrix} u_{0ij} \\ u_{1ij} \end{bmatrix} \sim N(0, \Omega_u) : \Omega_u = \begin{bmatrix} \sigma_{u0}^2 & \\ & \sigma_{u1}^2 \end{bmatrix}$$
$$\begin{bmatrix} e_{0ij} \end{bmatrix} \sim N(0, \Omega_e) : \Omega_e = \begin{bmatrix} \sigma_{e0}^2 \end{bmatrix}$$

At the bottom of the window is a toolbar with the following buttons: Fonts, Subs, Name, +, -, Add Term, Estimates, Nonlinear, Help (with a question mark icon), and Clea.

Figure 3. Equation screen with estimates.

Equations

$$y_{patient, clinic} \sim N(XB, \Omega)$$

$$y_{patient, clinic} = \beta_{0\ patient, clinic} x_0 + \beta_{1\ clinic} x_{1\ patient, clinic}$$

$$\beta_{0\ patient, clinic} = -0.012(0.040) + u_{0\ clinic} + e_{0\ patient, clinic}$$

$$\beta_{1\ clinic} = 0.557(0.020) + u_{1\ clinic}$$

$$\begin{bmatrix} u_{0\ clinic} \\ u_{1\ clinic} \end{bmatrix} \sim N(0, \Omega_u) : \Omega_u = \begin{bmatrix} 0.090(0.018) \\ 0.018(0.007) & 0.015(0.004) \end{bmatrix}$$

$$\begin{bmatrix} e_{0\ patient, clinic} \end{bmatrix} \sim N(0, \Omega_e) : \Omega_e = \begin{bmatrix} 0.554(0.012) \end{bmatrix}$$

$-2 * \loglikelihood(IGLS) = 9316.870(4059\ of\ 4059\ cases\ in\ use)$

Fonts Subs Name + - Add Term Estimates Nonlinear Help Clear

Figure 4

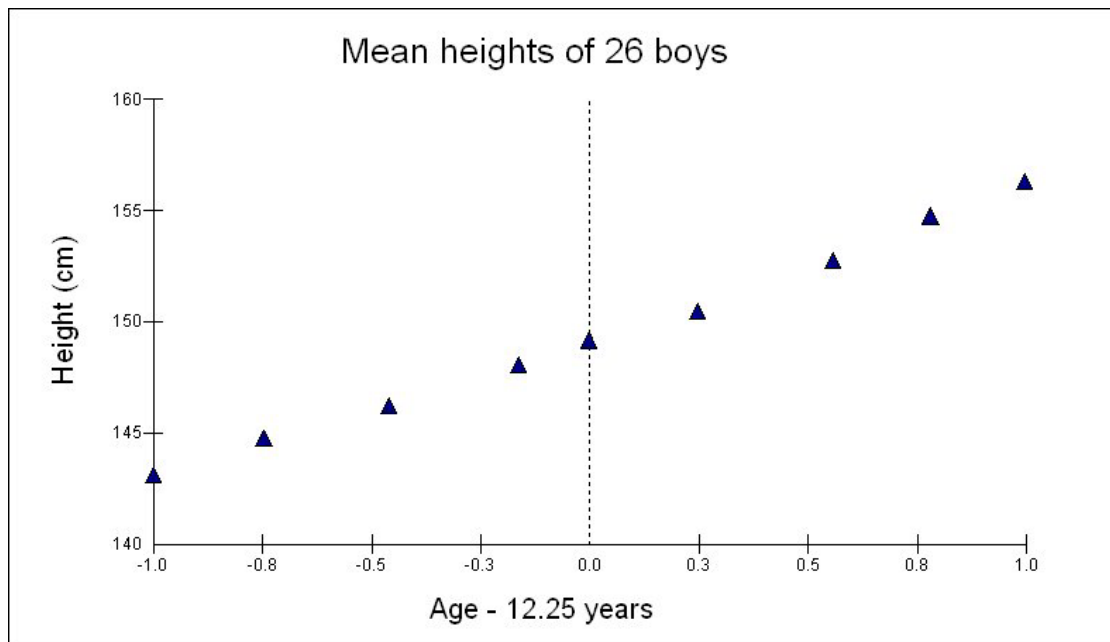


Figure 5

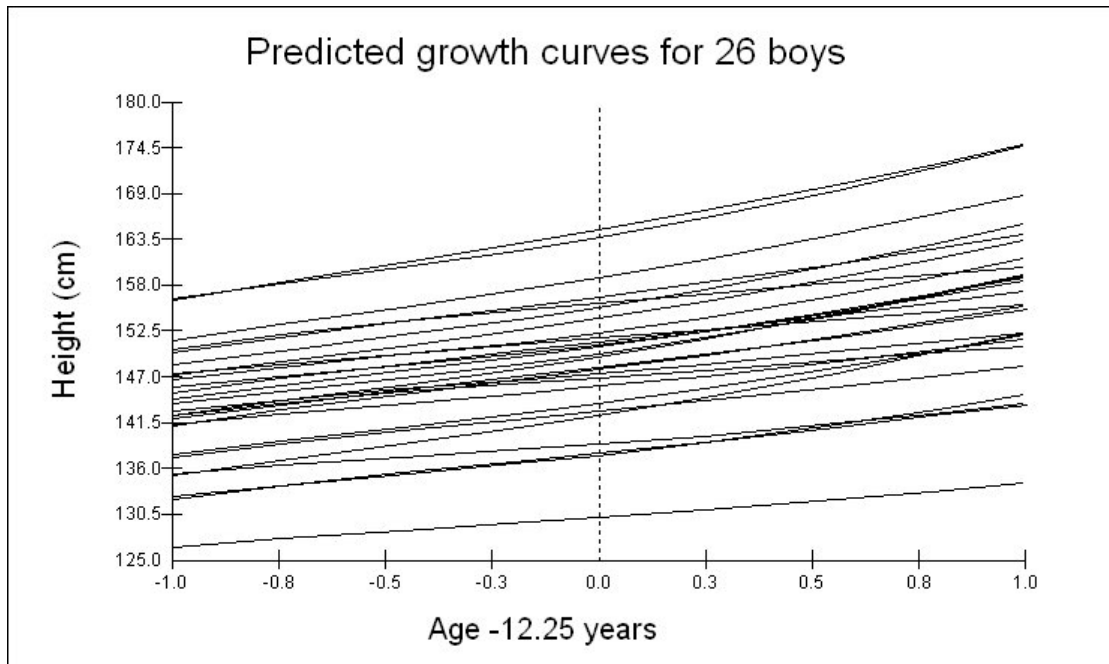


Figure 6. MCMC summary screen

