

Multilevel Analysis

Further Statistical Methods, Lecture 9 Hilary Term 2004

Steffen Lauritzen, University of Oxford; April 11, 2005

Multilevel observations

Multilevel analysis is concerned with observations with a *nested* structure.

For a two-level analysis we typically think of *individuals* within *groups*. The individual level is in general called *level one*, the group level *level two*.

An example of observations of this type can for example be performance measures for *pupils* of a specific age-group within *classes*.

The levels could be nested yet another time as e.g. classes within *schools*. And further, the schools could be grouped according to *regions* within *countries*, etc. although at the

top-level there might well be problems of compatibility of performance measures.

For simplicity we will only consider two levels, pupils within classes.

An example

As our basic example we will consider a Dutch study comprising $N = 131$ classes, each of sizes between 4 and 35, with a total of $N = 2287$ pupils.

The performance measure of interest is the score on a language test, and explanatory variables include class sizes and the IQ of individual pupils.

We let $Y_{ij}, j = 1, \dots, N, i = 1, \dots, n_j$ be the score for pupil j in class i and study the dependence of this response on covariates such as the IQ x_{ij} of the pupil and the size z_j of the class.

x_{ij} are *level one covariates* and z_j *level two covariates*.

A simple regression model

A first attempt could be to let

$$Y_{ij} = \beta_0 + \beta_1 x_{ij} + \beta_2 z_j + R_{ij}$$

with R_{ij} independent and distributed as $\mathcal{N}(0, \sigma^2)$.

This is a standard linear regression model which only has an indirect multilevel character.

The model ignores that pupils in the same class will tend to have more similar scores than those in different classes, even when the covariates are taken into account.

This is a *very serious mistake* if the variations in score at group level are not fully explained by the covariates.

Introducing random effects

For a moment, ignore the covariates x_{ij} and z_j and consider instead the model

$$Y_{ij} = \beta_0 + U_j + R_{ij}$$

where $U_j \sim \mathcal{N}(0, \tau^2)$. This model then has

$$\mathbf{V}(Y_{ij}) = \sigma^2 + \tau^2, \quad \text{Cov}(Y_{ij}, Y_{ij'}) = \tau^2, \quad \text{Cov}(Y_{ij}, Y_{i'j'}) = 0$$

so that scores of pupils within the same class are correlated. The correlation is

$$\rho = \frac{\tau^2}{\sigma^2 + \tau^2}$$

and is known as the *intraclass correlation coefficient*.

This type of model is also known as a *random effects model* since one could think of $\beta_j = \beta_0 + U_j$ as a group effect, in this case modelled as a random effect. Adding back the covariates leads to

$$Y_{ij} = \beta_0 + \beta_1 x_{ij} + \beta_2 z_j + U_j + R_{ij}.$$

It can give a better overview to introduce an intermediate variable describing the total class effect

$$M_j = \beta_0 + \beta_2 z_j + U_j; \quad Y_{ij} = M_j + \beta_1 x_{ij} + R_{ij}$$

where M_j now become missing data, or rather latent variables.

Estimation of parameters

The maximum likelihood (ML) estimates of the parameters can be obtained using the EM algorithm, treating M_j as missing variables.

For 'complete data', with M_j observed, the estimation problem splits into two simple linear regression problems

1. Estimating $(\beta_0, \beta_2, \tau^2)$ by regressing M_j on z_j ;
2. Estimating β_1, σ^2 by regressing $Y_{ij} - M_j$ on x_{ij}

Unfortunately the ML estimates of the variance components (σ^2, τ^2) can be very biased, as these do not

take into account the loss in degrees of freedom due to the estimation of regression coefficients.

Instead a method known as *residual maximum likelihood* or REML is often used.

An algorithm of EM type also exists for calculating the REML estimates, but this and other methods have also been implemented in generally available software.