

Lecture 9

Mixed models

Mixed models extend linear regression models in that additionally to the fixed coefficients (=fixed effects) also random effects are included (hence, the name: mixed effects models). Mixed models are used to model clustered, hierarchical or longitudinal data, where the random effects typically model some unobserved source of heterogeneity. For more details on these models consult Demidenko, E. (2013) *Mixed models: Theory and applications* and Jiang, J. (2007) *Linear and Generalized Linear Mixed Models and Their Applications*.

The data we consider are of the form Y_{ij} , $i = 1, \dots, n$, $j = 1, \dots, k_i$. That is, there are n subjects (clusters) and for each subject i there are k_i observations. Longitudinal data are observations of n subjects which are taken over time. A typical example of a longitudinal study is a measurement of some quantity of interest (height, weight, blood pressure, etc.) over time on n subjects. In clustered data the observations are typically taken at one time point, but the observations can be grouped by some feature.

A linear mixed model is written as

$$Y_i = X_i\beta + Z_i u_i + \epsilon_i, \quad i = 1, \dots, n,$$

where $Y_i, \epsilon_i \in \mathbb{R}^{k_i}$, $X_i \in \mathbb{R}^{k_i \times p}$, $Z_i \in \mathbb{R}^{k_i \times m}$, $\beta \in \mathbb{R}^p$ are fixed coefficients (effects) and $u_i \in \mathbb{R}^m$ are random variables with zero mean, which are independent on ϵ_i . That is, all subjects have the same fixed effect β and different random effects u_i . Fixed effects β are also called marginal or population average effects. The linear mixed model can also be represented as

$$Y = X\beta + Zu + \epsilon, \quad \mathbb{E} \begin{pmatrix} u \\ \epsilon \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \quad \text{cov} \begin{pmatrix} u \\ \epsilon \end{pmatrix} = \begin{pmatrix} \sigma_u^2 D & 0 \\ 0 & \sigma_\epsilon^2 \Sigma \end{pmatrix},$$

for $Y, \epsilon \in \mathbb{R}^N$, $N = \sum_{i=1}^n k_i$, $X \in \mathbb{R}^{N \times p}$, $Z = \text{blockdiag}(Z_1, \dots, Z_n) \in \mathbb{R}^{N \times nm}$, $u \in \mathbb{R}^{nm}$. Typically, a normality assumption is made both for ϵ and u .

Examples

1. One-way ANOVA (analysis of variance) model is typically written as

$$Y_{ij} = \beta_i + \epsilon_{ij}, \quad i = 1, \dots, n, \quad j = 1, \dots, k_i.$$

for fixed β_i , which sometimes also represented as $\beta_i = \beta + u_i$, where both β and u_i are fixed. Assuming that u_i are random, for example, normally distributed, leads to a model also known as variance components model, which is a mixed model:

$$Y_{ij} = \beta + u_i + \epsilon_{ij}, \quad \epsilon_{ij} \sim \mathcal{N}(0, \sigma^2), \quad u_i \sim \mathcal{N}(0, \sigma_u^2), \quad i = 1, \dots, n, \quad j = 1, \dots, k_i.$$

We can represent this model in the matrix notation as above setting $X = 1_N$, $Z = \text{blockdiag}(1_{k_1}, \dots, 1_{k_n})$.

2. Growth curve models typically have the form

$$Y_{ij} = X_i(\beta + u_i) + \epsilon_{ij}, \quad \epsilon_{ij} \sim \mathcal{N}(0, \sigma^2), \quad u_i \sim \mathcal{N}_p(0_p, \sigma_u^2 I_p), \quad i = 1, \dots, n, \quad j = 1, \dots, k_i,$$

where $X_i \in \mathbb{R}^{n \times p}$ represents the growth pattern. For example, if Y_{ij} is a height of an i th child at time point t_j , then one could take $X_i = (1, t_i, t_i^2)$.

3. In small area estimation is assumed that the means in “small areas” differ for a random amount:

$$Y_i = X_i \beta + 1_{n_i} u_i + \epsilon_i, \quad \epsilon_i \sim \mathcal{N}_{k_i}(0_{k_i}, \sigma^2 I_{k_i}), \quad u_i \sim \mathcal{N}(0, \sigma_u^2), \quad i = 1, \dots, n,$$

where in each i th small area there are k_i observations available.

Estimation in mixed models can be carried out using the maximum likelihood approach, if distributional assumptions on ϵ and u are made. To make the model identifiable we assume that $\sum_{i=1}^n X_i^t X_i$ is non-singular and that $N = \sum_{i=1}^n k_i > p$. Moreover, at least one matrix $Z_i^t Z_i$ is positive definite and $\sum_{i=1}^n (k_i - m) = N - nm > 0$.

Under normality assumptions $\epsilon \sim \mathcal{N}_N(0_N, \sigma_\epsilon^2 \Sigma)$ and $u_i \sim \mathcal{N}_m(0_m, \sigma_u^2 \tilde{D})$ we can represent the mixed model as

$$Y \sim \mathcal{N}_N(X\beta, \sigma_\epsilon^2 \Sigma + \sigma_u^2 Z D Z^t),$$

where $D = \text{blockdiag}(\tilde{D}, \dots, \tilde{D}) \in \mathbb{R}^{nm \times nm}$, so that the log-likelihood for the parameter vector $\theta = (\beta, u, \sigma_\epsilon^2, \sigma_u^2, \Sigma, D)$ becomes

$$\ell(\theta) = -\frac{1}{2} \left\{ N \log(2\pi\sigma_\epsilon^2) + \log |V| + \frac{(Y - X\beta)^t V^{-1} (Y - X\beta)}{\sigma_\epsilon^2} \right\},$$

where $V = \Sigma + \sigma_u^2 Z D Z^t / \sigma_\epsilon^2$. It is often assumed that $\Sigma = I_N$. Assumption that $\tilde{D} = I_m$ is not justified if the observations for each subject (=within a cluster) are taken over time, but is reasonable in many other settings. If $\tilde{D} \neq I_m$, then some parametric covariance structure is assumed (e.g. an autoregressive process), so that estimation of \tilde{D} is reduced to estimation of several parameters. In the following we consider a simple case $\tilde{D} = I_m$

and $\Sigma = I_N$, so that $V = V(\lambda) = I_N + ZZ^t/\lambda$ for $\lambda = \sigma_\epsilon^2/\sigma_u^2$.

It is easy to see that the estimator for β is given by

$$\hat{\beta} = (X^t V^{-1} X)^{-1} X^t V^{-1} Y,$$

which is a generalised least squares estimator. Note that this estimator depends on the unknown λ . Estimator for σ_ϵ^2 results in

$$\hat{\sigma}_\epsilon^2 = \frac{(Y - X\hat{\beta})^t V^{-1} (Y - X\hat{\beta})}{N} = \frac{Y^t V^{-1} (Y - X\hat{\beta})}{N}.$$

It remains to estimate λ . Using

$$\begin{aligned} \frac{\partial \log(|V|)}{\partial \lambda} &= \text{tr} \left(V^{-1} \frac{\partial V}{\partial \lambda} \right), \quad \frac{\partial V}{\partial \lambda} = -\frac{1}{\lambda^2} ZZ^t \\ \frac{\partial V^{-1}}{\partial \lambda} &= -V^{-1} \frac{\partial V}{\partial \lambda} V^{-1} \\ V^{-1} &= I_N - Z(Z^t Z + \lambda I_{nm})^{-1} Z^t \end{aligned}$$

we get the estimating equation from the profile likelihood

$$\begin{aligned} \frac{-2\partial \ell(\lambda; \hat{\beta}, \hat{\sigma}_\epsilon^2)}{\partial \lambda} &= \frac{1}{\hat{\sigma}_\epsilon^2} \frac{\partial \hat{\sigma}_\epsilon^2}{\partial \lambda} - \frac{1}{\lambda^2} \text{tr} [\{I_N - Z(Z^t Z + \lambda I_{nm})^{-1} Z^t\} ZZ^t] \\ &= \frac{1}{\lambda \hat{\sigma}_\epsilon^2} Y^t S (I_N - S) Y - \frac{1}{\lambda} \text{tr} \{(Z^t Z + \lambda I_{nm})^{-1} Z^t Z\}, \end{aligned}$$

where $S = C(C^t C + \lambda J)^{-1} C^t$ for $C = (X, Z)$ and $J = \text{diag}(0_p, 1_{nm})$. Finding a solution to

$$\frac{\partial \ell(\lambda; \hat{\beta}, \hat{\sigma}_\epsilon^2)}{\partial \lambda} = 0$$

gives an estimator for λ and hence σ_ϵ^2 and β .

It is easy to see that the maximum likelihood estimators for σ_ϵ^2 is biased. This in turn leads to inferior estimators of λ and β . It has been suggested in Patterson and Thompson (1971) and later in Harville (1974) to use the so-called *restricted* (also called *residual*) maximum likelihood. This approach is best justified in the Bayesian framework putting a flat non-informative prior on β . The restricted likelihood is defined to be

$$\ell_R(\theta) = \ell(\theta) - \frac{1}{2} \log |\sigma_\epsilon^2 X^t V^{-1} X|.$$

The restricted maximum likelihood (REML) estimator for β (with known λ) is unchanged, the estimator for σ_ϵ^2 becomes

$$\hat{\sigma}_\epsilon^2 = \frac{Y^t V^{-1} (Y - X\hat{\beta})}{N - p}$$

and λ is found by solving

$$\frac{-2\partial\ell(\lambda; \hat{\beta}, \hat{\sigma}_\epsilon^2)}{\partial\lambda} = \frac{1}{\lambda\hat{\sigma}_\epsilon^2} Y^t S (I_N - S) Y - \frac{1}{\lambda} \{\text{tr}(S) - p\} = 0.$$

Once λ is obtained, $\hat{\beta}$ follows immediately and finally an estimator for $\hat{\sigma}_\epsilon^2$ is calculated. Conditions for the existence of (restricted) maximum likelihood estimators in linear mixed models are given Demidenko (2013). Under mild regularity conditions all asymptotic properties of maximum likelihood estimators (consistency, asymptotic normality) apply.

We have estimated fixed parameters in the model and herewith the marginal effect β . However, in practice the conditional effects are also of interest. Since u is assumed to be random, one typically uses the term “predictor” for \hat{u} , which is obtained as a *best linear unbiased predictor* (BLUP). Henderson (1950) suggested to maximise the joint density (up to a constant) of u and Y in order to get estimators for β and u :

$$\left| \begin{pmatrix} \sigma_u^2 D & 0 \\ 0 & \sigma_\epsilon^2 \Sigma \end{pmatrix} \right|^{-1/2} \exp \left\{ -\frac{1}{2} \begin{pmatrix} u \\ \epsilon \end{pmatrix}^t \begin{pmatrix} \sigma_u^2 D & 0 \\ 0 & \sigma_\epsilon^2 \Sigma \end{pmatrix}^{-1} \begin{pmatrix} u \\ \epsilon \end{pmatrix} \right\},$$

where $\epsilon = Y - X\beta - Zu$. This leads to the criterion (up to terms independent of β and u)

$$(Y - X\beta - Zu)^t (\sigma_\epsilon^2 \Sigma)^{-1} (Y - X\beta - Zu) + u^t (\sigma_u^2 D)^{-1} u.$$

Note that this expression highlights that u are penalised. Setting as before $\Sigma = I_N$, $D = I_{nm}$ and $\lambda = \sigma_\epsilon^2 / \sigma_u^2$ we get the criterion

$$(Y - X\beta - Zu)^t (Y - X\beta - Zu) + \lambda u^t u = (Y - C\alpha)^t (Y - C\alpha) + \lambda \alpha^t J \alpha,$$

for $\alpha = (\beta, u)$. From this criterion we obtain for a known λ

$$\begin{aligned} \hat{\alpha} &= (C^t C + \lambda J)^{-1} C^t Y \\ \hat{\beta} &= (X^t V^{-1} X)^{-1} X^t Y \\ \hat{u} &= \frac{1}{\lambda} Z^t V^{-1} (Y - X\hat{\beta}). \end{aligned}$$

The predictor \hat{u} is also the best linear unbiased predictor of u (for known β and λ), which can be seen as follows. Assume we would like to obtain the best linear predictor of u from $Y - X\beta = Zu + \epsilon$. That is, we aim to find

$$\arg \min_{a, B} E \|Zu - \{a + B(Y - X\beta)\}\|^2,$$

which are

$$\begin{aligned} \hat{a} &= E(Zu) - E\{B(Y - X\beta)\} = 0 \\ \hat{B} &= \text{cov}(Zu, Y - X\beta) \text{cov}(Y - X\beta)^{-1} = Z \text{cov}(u) Z^t \text{cov}(Y)^{-1} = \sigma_u^2 Z Z^t V^{-1} / \sigma_\epsilon^2, \end{aligned}$$

leading to

$$\hat{u} = \frac{1}{\lambda} Z^t V^{-1} (Y - X\beta).$$

Once β and λ are replaced by estimators, \hat{u} is, strictly speaking, not BLUP.

Inference in mixed models is based on asymptotic normality.