

PANEL DATA ANALYSIS: LINEAR MIXED-EFFECT MODELS

datascience@berkeley

Theory

- For hierarchical data with a single level of grouping, we can formulate the classical LMM at a given level of a grouping factor as follows

$$\mathbf{y}_i = \mathbf{X}_i \boldsymbol{\beta} + \mathbf{Z}_i \mathbf{b}_i + \epsilon_i$$

where \mathbf{y}_i , \mathbf{X}_i , $\boldsymbol{\beta}$, and ϵ_i are the vectors of continuous responses, the design matrix, and the vector of residual errors for group i , $i = 1, 2, \dots, n$, and \mathbf{Z}_i and \mathbf{b}_i are the matrix of covariates and the corresponding vector of random effects:

$$\mathbf{Z}_i \equiv \begin{pmatrix} z_{i1}^{(1)} & z_{i1}^{(2)} & \cdots & z_{i1}^{(q)} \\ \vdots & \vdots & \ddots & \vdots \\ z_{in_i}^{(1)} & z_{in_i}^{(2)} & \cdots & z_{in_i}^{(q)} \end{pmatrix}$$

$z_i^{(k)}$

where each of the columns in \mathbf{Z}_i can be more compactly represented as $(\mathbf{z}^{(k)})_i$ and $\mathbf{b}_i \equiv (b_{i1}, \dots, b_{iq})^T$

- Similar to the design matrix \mathbf{X}_i , the matrix \mathbf{Z}_i contains known values of q covariates, with corresponding unobservable effects \mathbf{b}_i

Importantly, $\mathbf{b}_i \sim N_q(\mathbf{0}, \mathbf{D})$, $\epsilon_i \sim N_{n_i}(\mathbf{0}, \mathbf{R}_i)$ with $\mathbf{b}_i \perp \epsilon_i$

- That is, the residual errors ϵ_i for the same group are independent of the random effects \mathbf{b}_i .
- Also assume that vectors of random effects and residual errors for different groups are independent of each other: $\mathbf{b}_i \perp \epsilon_{i'}$ for $i \neq i'$

$$\mathcal{D} = \sigma^2 \mathbf{D} \quad \text{and} \quad \mathcal{R}_i = \sigma^2 \mathbf{R}_i$$

where σ^2 is an unknown scale parameter. In general, we will assume that \mathbf{D} and \mathbf{R}_i are positive-definite, unless stated otherwise


- In addition to the fixed-effects parameters \mathbf{b} for the covariates used in constructing the design matrix \mathbf{X}_i , the proposed model includes two random components: the within-group residual errors ϵ and the random effects \mathbf{b}_i for the covariates included in the matrix \mathbf{Z}_i . The presence of fixed and random effects of known variables gives rise to the name of this class of models.
- Some authors, researchers, partitioners refer to the proposed framework as a *two-level* or *two-stage* model. Following authors such as Pinheiro and Bates (2000), we will refer it to as a single-level linear mixed model.
- The framework established so far can easily be extended to multilevel grouped data.
- For instance, a model for data with two levels of grouping, with observations grouped into N first-level groups (indexed by $i = 1, \dots, N$), each with n_i second-level (sub-)groups (indexed by $j = 1, \dots, n_i$) containing n_{ij} observations, can be written as

$$\mathbf{y}_{ij} = \mathbf{X}_{ij}\boldsymbol{\beta} + \mathbf{Z}_{1,ij}\mathbf{b}_i + \mathbf{Z}_{2,ij}\mathbf{b}_{ij} + \boldsymbol{\epsilon}_{ij}$$

$$\mathbf{b}_i \sim \mathcal{N}_{q_1}(\mathbf{0}, \mathcal{D}_1), \quad \mathbf{b}_{ij} \sim \mathcal{N}_{q_2}(\mathbf{0}, \mathcal{D}_2), \quad \text{and} \quad \boldsymbol{\epsilon}_{ij} \sim \mathcal{N}_{n_{ij}}(\mathbf{0}, \mathcal{R}_{ij})$$

where the random vectors \mathbf{b}_i , \mathbf{b}_{ij} , and $\boldsymbol{\epsilon}_{ij}$ are independent of each other.

- In this setup, \mathbf{b}_i are the random effects associated with the first-level groups, while \mathbf{b}_{ij} are the random effects, independent of the first-level random effects, associated with the second-level groups. The design matrices $\mathbf{Z}_{1,ij}$ and $\mathbf{Z}_{2,ij}$ could be identical. Following *Pinheiro and Bates (2000)*, we refer this model as a *two-level LMM*.
- It follows that, for the classical LMMs, the conditional distribution, $f_{y|b}(\mathbf{y}_i | \mathbf{b}_i)$, of \mathbf{y}_i given \mathbf{b}_i is multivariate normal, with the mean and variance defined as:

$$\begin{aligned} E(\mathbf{y}_i | \mathbf{b}_i) &\equiv \boldsymbol{\mu}_i = \mathbf{X}_i \boldsymbol{\beta} + \mathbf{Z}_i \mathbf{b}_i \\ \text{Var}(\mathbf{y}_i | \mathbf{b}_i) &= \sigma^2 \mathbf{R}_i, \end{aligned}$$


with $\boldsymbol{\mu}_i \equiv (\mu_{i1}, \dots, \mu_{i,n_i})'$

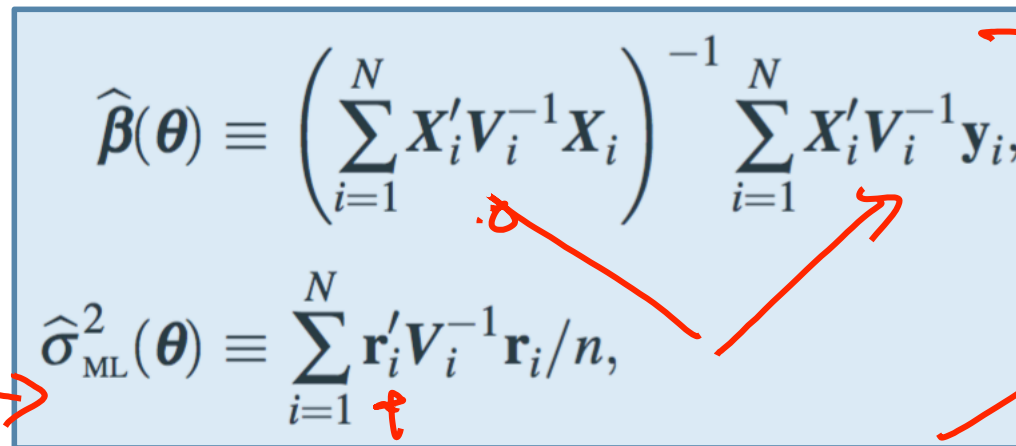
$$E(y_{ij} | \mathbf{b}_i) \equiv \mu_{ij} = \mathbf{x}_{ij}' \boldsymbol{\beta} + \mathbf{z}_{ij}' \mathbf{b}_i$$

where $\mathbf{x}_{ij} \equiv (x_{ij}^{(1)}, \dots, x_{ij}^{(p)})'$ and $\mathbf{z}_{ij} \equiv (z_{ij}^{(1)}, \dots, z_{ij}^{(q)})'$ are column vectors

- Conditionally on the (unknown) values of the random effects \mathbf{b}_i , the mean value of the dependent-variable vector \mathbf{y}_i is defined by a linear combination of the vectors of the \mathbf{X} - and \mathbf{Z} -covariates included, as columns, in the group-specific design matrices \mathbf{X}_i and \mathbf{Z}_i , corresponding to the fixed effects β and random effects \mathbf{b}_i .
- Moreover, the conditional variance-covariance matrix of \mathbf{y}_i is equal to the variance-covariance matrix of the residual errors ϵ_i

Maximum-Likelihood Estimation

$$\begin{aligned}\ell_{\text{Full}}(\boldsymbol{\beta}, \sigma^2, \boldsymbol{\theta}) &\equiv -\frac{N}{2} \log(\sigma^2) - \frac{1}{2} \sum_{i=1}^N \log[\det(\mathbf{V}_i)] \\ &\quad - \frac{1}{2\sigma^2} \sum_{i=1}^N (\mathbf{y}_i - \mathbf{X}_i \boldsymbol{\beta})' \mathbf{V}_i^{-1} (\mathbf{y}_i - \mathbf{X}_i \boldsymbol{\beta})\end{aligned}$$



$$\begin{aligned}\hat{\boldsymbol{\beta}}(\boldsymbol{\theta}) &\equiv \left(\sum_{i=1}^N \mathbf{X}_i' \mathbf{V}_i^{-1} \mathbf{X}_i \right)^{-1} \sum_{i=1}^N \mathbf{X}_i' \mathbf{V}_i^{-1} \mathbf{y}_i, \\ \hat{\sigma}_{\text{ML}}^2(\boldsymbol{\theta}) &\equiv \sum_{i=1}^N \mathbf{r}_i' \mathbf{V}_i^{-1} \mathbf{r}_i / n,\end{aligned}$$

where $\mathbf{r}_i \equiv \mathbf{r}_i(\boldsymbol{\theta}) = \mathbf{y}_i - \mathbf{X}_i \hat{\boldsymbol{\beta}}(\boldsymbol{\theta})$.

$$\hat{\sigma}_{\text{REML}}^2(\boldsymbol{\theta}) \equiv \sum_{i=1}^N \mathbf{r}_i' \mathbf{V}_i^{-1} \mathbf{r}_i / (n - p)$$

- This leads to a log-profile-restricted-likelihood function, which only depends on θ :

$$\begin{aligned} \ell_{\text{REML}}^*(\boldsymbol{\theta}) \equiv & -\frac{n-p}{2} \log \left(\sum_{i=1}^N \mathbf{r}_i' \mathbf{r}_i \right) - \frac{1}{2} \sum_{i=1}^N \log[\det(\mathbf{V}_i)] \\ & - \frac{1}{2} \log \left[\det \left(\sum_{i=1}^N \mathbf{X}_i' \mathbf{V}_i^{-1} \mathbf{X}_i \right) \right]. \end{aligned}$$

To sum up, we have:

$$\mathcal{Y} \sim \mathcal{N}(\mathbf{X}\boldsymbol{\beta} + \mathbf{o}, \sigma^2 \mathbf{W}^{-1})$$

$$(\mathcal{Y}|\mathcal{B} = \mathbf{b}) \sim \mathcal{N}(\mathbf{X}\boldsymbol{\beta} + \mathbf{Zb} + \mathbf{o}, \sigma^2 \mathbf{W}^{-1})$$

$\mathcal{B} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma})$ a parameterized $q \times q$ variance-covariance matrix, $\boldsymbol{\Sigma}$.

As a variance-covariance matrix, $\boldsymbol{\Sigma}$ must be positive semidefinite.

$$\left(\boldsymbol{\Sigma}_{\boldsymbol{\theta}} = \sigma^2 \boldsymbol{\Lambda}_{\boldsymbol{\theta}} \boldsymbol{\Lambda}_{\boldsymbol{\theta}}^{\top}; \right.$$

a relative covariance factor, $\boldsymbol{\Lambda}_{\boldsymbol{\theta}}$, which is a $q \times q$ matrix,

Final Remarks :

- This formulation of linear mixed models allows for a relatively compact expression for the profiled log-likelihood of θ
- The matrices associated with random effects, \mathbf{Z} and Λ_θ , typically have a sparse structure with a sparsity pattern that encodes various model assumptions.
- More on the details on the structure and how to represent it efficiently can be found in the assigned readings, R demo, or the live sessions.
- The setup above is very general. The interface provided by *lme4*'s *lmer* function is less general than the model described.
- To take advantage of the entire range of possibilities, one may have to use the modular functions or explore the experimental *flexLambda* branch of *lme4* on *Github*

Berkeley

SCHOOL OF
INFORMATION