

STAT 574 Linear and Nonlinear Mixed Models

Lecture 2: Linear Mixed Models and Estimation

Chencheng Cai

Washington State University

Weight versus height example

- ▶ A dataset contains the heights and weights of 71 people from 18 families.
- ▶ A naive linear regression model:

$$W_k = \alpha + \beta H_k + \epsilon_k,$$

where W_k is the weight of the k th person and H_k is his/her height.

Anything wrong with the linear regression assumptions?

LINE Assumptions

- ▶ Linear: it might be more realistic to assume

$$W_k = \alpha + \beta H_k^2 + \epsilon_k.$$

Justification: Body Mass Index (BMI)

- ▶ Independence: weights of people in the same family are highly correlated.
Potential confounders: gene, habits, environment.
- ▶ Normal: can be diagnosed later.
- ▶ Equal variance: may need to take a logarithm of the weight.

LINE Assumptions

- ▶ Linear: it might be more realistic to assume

$$W_k = \alpha + \beta H_k^2 + \epsilon_k.$$

Justification: Body Mass Index (BMI)

- ▶ Independence: weights of people in the same family are highly correlated.
Potential confounders: gene, habits, environment.
- ▶ Normal: can be diagnosed later.
- ▶ Equal variance: may need to take a logarithm of the weight.

let us focus on the correlation for now.

Model update

In order to incorporate the within-family correlation, we assume

$$W_{ij} = \alpha_i + \beta H_{ij} + \epsilon_{ij}$$

- ▶ The index ij refers to the j th person in the i th family.
- ▶ Intercept α_i is family-specific.
- ▶ ϵ_{ij} is normal-distributed and is independent with other variables.
- ▶ α_i 's are I.I.D.

Correlation

- ▶ Alice and Bob from the same family:

$$\begin{aligned}\text{Cov}(W_{ij}, W_{ij'}) &= \text{Cov}(\alpha_i + \beta H_{ij} + \epsilon_{ij}, \alpha_i + \beta H_{ij'} + \epsilon_{ij'}) \\ &= \text{Cov}(\alpha_i + \epsilon_{ij}, \alpha_i + \epsilon_{ij'}) \\ &= \text{Var}(\alpha_i) > 0\end{aligned}$$

- ▶ Alice and Bob from different families:

$$\begin{aligned}\text{Cov}(W_{ij}, W_{i'j'}) &= \text{Cov}(\alpha_i + \beta H_{ij} + \epsilon_{ij}, \alpha_{i'} + \beta H_{i'j'} + \epsilon_{i'j'}) \\ &= \text{Cov}(\alpha_i + \epsilon_{ij}, \alpha_{i'} + \epsilon_{i'j'}) \\ &= \text{Cov}(\alpha_i, \alpha_{i'}) \\ &= 0\end{aligned}$$

Model reformat

We can also write

$$W_{ij} = \alpha + \beta H_{ij} + b_i \times 1 + \epsilon_{ij}$$

where $\alpha = \mathbb{E}[\alpha_i]$ and $b_i = \alpha_i - \alpha$.

Furthermore, we vectorize it by family:

$$\mathbf{W}_i = \alpha + \beta \mathbf{H}_i + b_i \mathbf{Z}_i + \boldsymbol{\epsilon}_i$$

where $\mathbf{Z}_i = (1, 1, \dots, 1)^T$.

Linear Mixed Effects (LME) Model

$$\mathbf{y}_i = \mathbf{X}_i\boldsymbol{\beta} + \mathbf{Z}_i\mathbf{b}_i + \boldsymbol{\epsilon}_i, \quad \text{for } i = 1, \dots, N.$$

- ▶ \mathbf{y}_i : $n_i \times 1$ vector of responses of the i th cluster/group.
- ▶ \mathbf{X}_i : $n_i \times m$ design matrix of fixed effects.
- ▶ $\boldsymbol{\beta}$: $m \times 1$ vector of fixed effects coefficients.
- ▶ \mathbf{Z}_i : $n_i \times k$ design matrix of random effects.
- ▶ \mathbf{b}_i : $k \times 1$ vector of random effect coefficients.
- ▶ $\boldsymbol{\epsilon}_i$: $n_i \times 1$ vector of error terms.

Examples

- ▶ Weight-height relation for different families.
- ▶ Social behaviors for people of different ages.
- ▶ Stock price movement for companies in different sections/industries.
- ▶ Biological studies on animals of different sub-species.

In summary, you should consider a mixed-effects model if you do not believe your model (and the coefficients) are the same for different sub-populations.

Distributional Assumptions

Normal assumptions for the error terms and for the random effect coefficients:

$$\epsilon_i \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}), \quad \mathbf{b}_i \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{D})$$

warning! Sometimes, one only assumes $\text{Cov}(\mathbf{b}_i) = \sigma^2 \mathbf{D}$. But we assume normal here.

Consequently,

$$\mathbf{y}_i \sim \mathcal{N}(\mathbf{X}_i \boldsymbol{\beta}, \sigma^2 \mathbf{V}_i), \quad \text{for } i = 1, \dots, N$$

with $\mathbf{V}_i = \mathbf{I} + \mathbf{Z}_i \mathbf{D} \mathbf{Z}_i^T$.

A more compact formula

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\eta}$$

where

$$\mathbf{y} = \begin{bmatrix} \mathbf{y}_1 \\ \mathbf{y}_2 \\ \vdots \\ \mathbf{y}_N \end{bmatrix}, \quad \mathbf{X} = \begin{bmatrix} \mathbf{X}_1 \\ \mathbf{X}_2 \\ \vdots \\ \mathbf{X}_N \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} \mathbf{b}_1 \\ \mathbf{b}_2 \\ \vdots \\ \mathbf{b}_N \end{bmatrix}, \quad \boldsymbol{\epsilon} = \begin{bmatrix} \boldsymbol{\epsilon}_1 \\ \boldsymbol{\epsilon}_2 \\ \vdots \\ \boldsymbol{\epsilon}_N \end{bmatrix}$$

and $\mathbf{A} = \text{diag}(\mathbf{Z}_1, \mathbf{Z}_2, \dots, \mathbf{Z}_N)$. Also

$$\text{Cov}(\boldsymbol{\eta}) = \sigma^2 \text{diag}(\mathbf{V}_1, \dots, \mathbf{V}_N)$$

Special Model: random intercepts model

$$y_{ij} = \alpha_i + \boldsymbol{\gamma}'\mathbf{u}_{ij} + \epsilon_{ij}.$$

where $\alpha_i = \alpha + b_i$ and $b_i \sim \mathcal{N}(0, \sigma^2 d)$.

Special Model: balanced random-coefficient model

We assume all clusters have the same size $n_i = n$ and

$$\mathbf{Z} = \mathbf{X}_i = \mathbf{Z}_i, \quad i = 1, \dots, N.$$

Then, LME becomes:

$$\mathbf{y}_i = \mathbf{Z}\boldsymbol{\beta} + \mathbf{Z}\mathbf{b}_i + \boldsymbol{\epsilon}_i, \quad i = 1, \dots, N.$$

A more compact form:

$$\mathbf{Y} = \mathbf{Z}\boldsymbol{\beta}\mathbf{1}^T + \mathbf{E}.$$

with $\mathbf{Y} = [\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_N]$ and \mathbf{E} has i.i.d. columns with covariance $\sigma^2(\mathbf{I} + \mathbf{Z}\mathbf{D}\mathbf{Z}^T)$.

Linear Growth Curve Model

We start with a linear regression model with random coefficients:

$$\mathbf{y}_i = \mathbf{Z}_i \mathbf{a}_i + \epsilon_i.$$

Furthermore, we assume the coefficients are linear combinations of other covariates:

$$\mathbf{a}_i = \mathbf{A}_i \boldsymbol{\beta} + \mathbf{b}_i$$

with $\mathbb{E}(\mathbf{b}_i) = 0$ and $\text{Cov}(\mathbf{b}_i) = \sigma^2 \mathbf{D}$.

Combine them:

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

Example:

- ▶ Index i : person. index j : time.
- ▶ y : health indicator. \mathbf{Z} : health-related covariates.
- ▶ \mathbf{A} : whether the person takes a medicine or a placebo. $\boldsymbol{\beta}$: the effect of the medicine.

Log-likelihood Function

- ▶ Distribution:

$$\mathbf{y}_i \sim \mathcal{N}(\mathbf{X}_i\boldsymbol{\beta}, \sigma^2\mathbf{V}_i), \quad \mathbf{V}_i = \mathbf{I} + \mathbf{Z}_i\mathbf{D}\mathbf{Z}_i^T.$$

- ▶ Probability:

$$p(\mathbf{y}_i \mid \mathbf{X}_i, \mathbf{Z}_i, \boldsymbol{\beta}, \mathbf{D}, \sigma^2) = \frac{1}{(2\pi)^{n_i/2} \sqrt{|\sigma^2\mathbf{V}_i|}} \exp \left\{ -\frac{1}{2\sigma^2} (\mathbf{y}_i - \mathbf{X}_i\boldsymbol{\beta})^T \mathbf{V}_i^{-1} (\mathbf{y}_i - \mathbf{X}_i\boldsymbol{\beta}) \right\}$$

- ▶ Log-likelihood:

$$\begin{aligned} \ell(\boldsymbol{\beta}, \mathbf{D}, \sigma^2) &= \sum_{i=1}^N \left\{ -\frac{n_i}{2} \log 2\pi\sigma^2 - \frac{1}{2} \log |\mathbf{V}_i| - \frac{1}{2\sigma^2} (\mathbf{y}_i - \mathbf{X}_i\boldsymbol{\beta})^T \mathbf{V}_i^{-1} (\mathbf{y}_i - \mathbf{X}_i\boldsymbol{\beta}) \right\} \\ &= -\frac{1}{2} \left\{ N_T \log \sigma^2 + \sum_{i=1}^N \left[\log |\mathbf{V}_i| + \frac{1}{\sigma^2} (\mathbf{y}_i - \mathbf{X}_i\boldsymbol{\beta})^T \mathbf{V}_i^{-1} (\mathbf{y}_i - \mathbf{X}_i\boldsymbol{\beta}) \right] \right\} + C \end{aligned}$$

where $N_T = \sum_{i=1}^N n_i$ and $C = -\frac{N_T}{2} \log 2\pi$.

MLE for LME Models

$$\max_{(\boldsymbol{\beta}, \mathbf{D}, \sigma^2) \in \Theta} \underbrace{-\frac{1}{2} \left\{ N_T \log \sigma^2 + \sum_{i=1}^N [\log |\mathbf{V}_i| + \sigma^{-2} \mathbf{e}_i^T \mathbf{V}_i^{-1} \mathbf{e}_i] \right\}}_{\ell(\boldsymbol{\beta}, \mathbf{D}, \sigma^2)}$$

with $\mathbf{V}_i = \mathbf{I} + \mathbf{Z}_i \mathbf{D} \mathbf{Z}_i^T$ and $\mathbf{e}_i = \mathbf{y}_i - \mathbf{X}_i \boldsymbol{\beta}$.

Optimization algorithms:

- ▶ Newton–Raphson — optimized gradient descent.
- ▶ Fisher scoring — a stable version of NR
- ▶ Expectation–Maximization — for missing value problems

The Parameter Space

- ▶ Nonnegative Definite Parameter Space:

$$\Theta = \{(\boldsymbol{\beta}, \mathbf{D}, \sigma^2) : \boldsymbol{\beta} \in \mathbb{R}^m, \sigma^2 > 0, \mathbf{D} \succeq 0\}$$

- ▶ Dimension: $\dim(\Theta) = m + 1 + k(k + 1)/2$.
- ▶ Drawback: difficult to enforce nonnegativeness.
- ▶ Another parameter space:

$$\Theta = \{(\boldsymbol{\beta}, \mathbf{D}, \sigma^2) : \boldsymbol{\beta} \in \mathbb{R}^m, \sigma^2 > 0, \mathbf{V}_i \succ 0 \text{ for } i = 1, \dots, N\}$$

- ▶ Benefits: $\ell(\theta) \rightarrow -\infty$ on boundary.

Profiling LLH Function

Consider the one-step optimization:

$$\max_{\sigma^2} -\frac{1}{2} \left\{ N_T \log \sigma^2 + \sum_{i=1}^N [\log |\mathbf{V}_i| + \sigma^{-2} \mathbf{e}_i^T \mathbf{V}_i^{-1} \mathbf{e}_i] \right\}$$

Partial derivative:

$$\frac{\partial \ell}{\partial \sigma^2} = -\frac{N_T}{2\sigma^2} + \frac{1}{2\sigma^4} \sum_{i=1}^N \mathbf{e}_i^T \mathbf{V}_i^{-1} \mathbf{e}_i.$$

Setting above to zero, we have:

$$\hat{\sigma}^2 = \frac{1}{N_T} \sum_{i=1}^N \mathbf{e}_i^T \mathbf{V}_i^{-1} \mathbf{e}_i$$

profiled log-likelihood function:

$$\ell_p(\boldsymbol{\beta}, \mathbf{D}) = \ell(\boldsymbol{\beta}, \mathbf{D}, \hat{\sigma}^2) = -\frac{1}{2} \left\{ N_T \log \sum_{i=1}^N \mathbf{e}_i^T \mathbf{V}_i^{-1} \mathbf{e}_i + \sum_{i=1}^N \log |\mathbf{V}_i| \right\} + C$$

Profiling LLH Function

One step further:

$$\max_{\boldsymbol{\beta}} \underbrace{-\frac{1}{2} \left\{ N_T \log \sum_{i=1}^N \mathbf{e}_i^T \mathbf{V}_i^{-1} \mathbf{e}_i + \sum_{i=1}^N \log |\mathbf{V}_i| \right\}}_{\ell_p(\boldsymbol{\beta}, \mathbf{D})}$$

Equivalent to:

$$\min_{\boldsymbol{\beta}} \sum_{i=1}^N \mathbf{e}_i^T \mathbf{V}_i^{-1} \mathbf{e}_i$$

Set the partial derivative to zero:

$$2 \sum_{i=1}^N (\mathbf{y}_i - \mathbf{X}_i \boldsymbol{\beta})^T \mathbf{V}_i^{-1} \mathbf{X}_i = 0$$

Profiling LLH Function

Solution:

$$\hat{\boldsymbol{\beta}} = \left(\sum_{i=1}^N \mathbf{X}_i^T \mathbf{V}_i^{-1} \mathbf{X}_i \right)^{-1} \left(\sum_{i=1}^N \mathbf{X}_i^T \mathbf{V}_i^{-1} \mathbf{y}_i \right)$$

The above solution is the Generalized Least Squares (GLS) estimator.

Further profiled log-likelihood function:

$$\ell_p(\mathbf{D}) = \ell_p(\hat{\boldsymbol{\beta}}, \mathbf{D}) = -\frac{1}{2} \left\{ N_T \log (s_{yy} - \mathbf{s}_{xy}^T \mathbf{S}_{xx}^{-1} \mathbf{s}_{xy}) + \sum_{i=1}^N \log |\mathbf{V}_i| \right\}$$

where $s_{yy} = \sum_i \mathbf{y}_i^T \mathbf{V}_i^{-1} \mathbf{y}_i$, $\mathbf{s}_{xy} = \sum_i \mathbf{X}_i^T \mathbf{V}_i^{-1} \mathbf{y}_i$ and $\mathbf{S}_{xx} = \sum_i \mathbf{X}_i^T \mathbf{V}_i^{-1} \mathbf{X}_i$.

We can maximize $\ell_p(\mathbf{D})$ instead of $\ell(\boldsymbol{\beta}, \mathbf{D}, \sigma^2)$ to speed up computation!

Profiling LLH Function

The final optimization:

$$\max_{\mathbf{D}} -\frac{1}{2} \left\{ N_T \log (s_{yy} - \mathbf{s}_{xy}^T \mathbf{S}_{xx}^{-1} \mathbf{s}_{xy}) + \sum_{i=1}^N \log |\mathbf{V}_i| \right\}$$

Unfortunately, it usually has no analytical solution.

Let $\hat{\mathbf{D}}$ be the optimum and $\hat{\mathbf{V}}_i = \mathbf{I} + \mathbf{Z}_i \hat{\mathbf{D}} \mathbf{Z}_i^T$. Then, we have the MLE for the other two parameters:

$$\hat{\boldsymbol{\beta}} = \left(\sum_{i=1}^N \mathbf{X}_i^T \hat{\mathbf{V}}_i^{-1} \mathbf{X}_i \right)^{-1} \left(\sum_{i=1}^N \mathbf{X}_i^T \hat{\mathbf{V}}_i^{-1} \mathbf{y}_i \right)$$
$$\hat{\sigma}^2 = \frac{1}{N_t} \sum_{i=1}^N (\mathbf{y}_i - \mathbf{X}_i \hat{\boldsymbol{\beta}})^T \hat{\mathbf{V}}_i^{-1} (\mathbf{y}_i - \mathbf{X}_i \hat{\boldsymbol{\beta}})$$

Optimizing $\ell_p(\mathbf{D})$

Using Woodbury identity $V_i^{-1} = \mathbf{I} - \mathbf{Z}_i(\mathbf{D}^{-1} + \mathbf{Z}_i^T \mathbf{Z}_i)^{-1} \mathbf{Z}_i^T$, we have

$$\begin{aligned} s_{yy} &= \sum \mathbf{y}_i^T \mathbf{y}_i - \sum \mathbf{y}_i^T \mathbf{Z}_i (\mathbf{D}^{-1} + \mathbf{Z}_i^T \mathbf{Z}_i)^{-1} \mathbf{Z}_i^T \mathbf{y}_i \\ s_{xy} &= \sum \mathbf{X}_i^T \mathbf{y}_i - \sum \mathbf{X}_i^T \mathbf{Z}_i (\mathbf{D}^{-1} + \mathbf{Z}_i^T \mathbf{Z}_i)^{-1} \mathbf{Z}_i^T \mathbf{y}_i \\ s_{xx} &= \sum \mathbf{X}_i^T \mathbf{X}_i - \sum \mathbf{X}_i^T \mathbf{Z}_i (\mathbf{D}^{-1} + \mathbf{Z}_i^T \mathbf{Z}_i)^{-1} \mathbf{Z}_i^T \mathbf{X}_i \end{aligned}$$

Using Sylvester's identity $|\mathbf{I} + \mathbf{Z}_i \mathbf{D} \mathbf{Z}_i^T| = |\mathbf{I} + \mathbf{D} \mathbf{Z}_i^T \mathbf{Z}_i|$, we have

$$\log |\mathbf{V}_i| = \log(|\mathbf{D}| |\mathbf{D}^{-1} + \mathbf{Z}_i^T \mathbf{Z}_i|) = \log |\mathbf{D}^{-1} + \mathbf{Z}_i^T \mathbf{Z}_i| - \log |\mathbf{D}^{-1}|$$

Benefits:

- ▶ Complicated computation (determinant, inverse) on $n_i \times n_i$ matrices are now on $k \times k$ matrices.
- ▶ Many quantities can be pre-computed.
- ▶ Optimization can be done with respect to \mathbf{D}^{-1} .

Optimizing $\ell_p(\mathbf{D}^{-1})$

For algorithms, it is necessary to know the partial derivative of $\ell_p(\mathbf{D}^{-1})$.

$$d \log |\mathbf{V}_i| = \text{tr} [(\mathbf{D}^{-1} + \mathbf{Z}_i^T \mathbf{Z}_i)^{-1} d\mathbf{D}^{-1}] - \text{tr} [\mathbf{D} d\mathbf{D}^{-1}]$$

$$ds_{yy} = \sum \mathbf{y}_i^T \mathbf{Z}_i (\mathbf{D}^{-1} + \mathbf{Z}_i^T \mathbf{Z}_i)^{-1} d\mathbf{D}^{-1} (\mathbf{D}^{-1} + \mathbf{Z}_i^T \mathbf{Z}_i)^{-1} \mathbf{Z}_i^T \mathbf{y}_i$$

$$ds_{xy} = \sum \mathbf{x}_i^T \mathbf{Z}_i (\mathbf{D}^{-1} + \mathbf{Z}_i^T \mathbf{Z}_i)^{-1} d\mathbf{D}^{-1} (\mathbf{D}^{-1} + \mathbf{Z}_i^T \mathbf{Z}_i)^{-1} \mathbf{Z}_i^T \mathbf{y}_i$$

$$d\mathbf{S}_{xx} = \sum \mathbf{x}_i^T \mathbf{Z}_i (\mathbf{D}^{-1} + \mathbf{Z}_i^T \mathbf{Z}_i)^{-1} d\mathbf{D}^{-1} (\mathbf{D}^{-1} + \mathbf{Z}_i^T \mathbf{Z}_i)^{-1} \mathbf{Z}_i^T \mathbf{x}_i$$

Optimizing $\ell_p(\mathbf{D}^{-1})$

Eventually, we have

$$\frac{\partial \ell_p(\mathbf{D}^{-1})}{\partial \mathbf{D}^{-1}} = -\frac{1}{2} \left\{ \sum_{i=1}^N \mathbf{G}_i - N\mathbf{D} + \frac{N_T}{s_{yy} - \mathbf{s}_{xy}^T \mathbf{S}_{xx} \mathbf{s}_{xy}} \sum_{i=1}^N \mathbf{G}_i \mathbf{Z}_i^T \mathbf{e}_i \mathbf{e}_i^T \mathbf{Z}_i \mathbf{G}_i \right\}$$

where

$$\mathbf{G}_i = (\mathbf{D}_i^{-1} + \mathbf{Z}_i^T \mathbf{Z}_i)^{-1}$$

$$\mathbf{e}_i = \mathbf{y}_i - \mathbf{X}_i \mathbf{S}_{xx}^{-1} \mathbf{s}_{xy}$$

Restricted MLE

The MLE for σ^2 is estimated from $\hat{e}_i = \mathbf{y}_i - \mathbf{X}_i\hat{\boldsymbol{\beta}}$ directly:

$$\hat{\sigma}^2 = \frac{1}{N_t} \sum_{i=1}^N \hat{e}_i^T \hat{\mathbf{V}}_i^{-1} \hat{e}_i$$

$\hat{\sigma}^2$ is in general biased for σ^2 because \hat{e}_i involves $\hat{\boldsymbol{\beta}}$.

The likelihood function (after certain modifications) that disentangles $\hat{\boldsymbol{\beta}}$ from other estimators is called the **restricted likelihood** function.

The method that maximizes the restricted likelihood function is called **restricted maximum likelihood (REML)**.

REML

- ▶ Consider a linear regression model $\mathbf{y} \sim \mathcal{N}(\mathbf{X}\boldsymbol{\beta}, \mathbf{V})$.
- ▶ The observation can be decomposed into two orthogonal parts: $\mathbf{X}\hat{\boldsymbol{\beta}}$ and $\hat{\mathbf{e}} = \mathbf{y} - \mathbf{X}\boldsymbol{\beta}$.
- ▶ Why orthogonal?

$$\text{Cov}(\mathbf{X}\hat{\boldsymbol{\beta}}, \hat{\mathbf{e}}) = \text{Cov}(\mathbf{H}\mathbf{y}, (\mathbf{I} - \mathbf{H})\mathbf{y}) = \mathbf{H}\mathbf{V}(\mathbf{I} - \mathbf{H}) = \mathbf{0}.$$

- ▶ Therefore, we can write:

$$\ell(\hat{\mathbf{e}}, \mathbf{V}) = \ell(\mathbf{y}, \mathbf{V}) - \ell(\hat{\boldsymbol{\beta}}, \mathbf{V}) = -\frac{1}{2} \{ \log |\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X}| + \log |\mathbf{V}| + \hat{\mathbf{e}}^T \mathbf{V}^{-1} \hat{\mathbf{e}} \}$$

- ▶ Furthermore, we can have the restricted log-likelihood:

$$\ell_R(\boldsymbol{\beta}, \mathbf{V}) = -\frac{1}{2} \{ \log |\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X}| + \log |\mathbf{V}| + (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T \mathbf{V}^{-1} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \}$$

REML — Bayesian Perspective

We take a non-informative prior on β (i.e. uniform). Then the profiled likelihood function is

$$L_R(\mathbf{V}) = \int_{\mathbb{R}^m} L(\beta, \mathbf{V}) d\beta \propto |\mathbf{V}|^{-1/2} |\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X}|^{-1/2} \exp \left\{ -\frac{1}{2} (\mathbf{y} - \mathbf{X} \hat{\beta})^T \mathbf{V}^{-1} (\mathbf{y} - \mathbf{X} \hat{\beta}) \right\}$$

- ▶ Why non-informative?
- ▶ Why independent?

REML for LME

- ▶ Back to linear mixed models, we have

$$\ell_R(\boldsymbol{\beta}, \mathbf{D}, \sigma^2) = -\frac{1}{2} \left\{ (N_T - m) \log \sigma^2 + \log \left| \sum_{i=1}^N \mathbf{X}_i^T \mathbf{V}_i^{-1} \mathbf{X}_i \right| \right. \\ \left. + \sum_{i=1}^N [\log |\mathbf{V}_i| + \sigma^{-2} (\mathbf{y}_i - \mathbf{X}_i \boldsymbol{\beta})^T \mathbf{V}_i^{-1} (\mathbf{y}_i - \mathbf{X}_i \boldsymbol{\beta})] \right\}$$

- ▶ $\log \left| \sum_{i=1}^N \mathbf{X}_i^T \mathbf{V}_i^{-1} \mathbf{X}_i \right|$ v.s. $\sum_{i=1}^N \log |\mathbf{V}_i|$
- ▶ $(N_T - m) \log \sigma^2$

Profiled REML

- ▶ Maximize $\ell_R(\boldsymbol{\beta}, \mathbf{D}, \sigma^2)$:

$$\hat{\sigma}_R^2 = \frac{1}{N_T - m} \sum_{i=1}^N (\mathbf{y}_i - \mathbf{X}_i \boldsymbol{\beta})^T \mathbf{V}_i^{-1} (\mathbf{y}_i - \mathbf{X}_i \boldsymbol{\beta})$$

- ▶ Therefore the profiled restricted LLH is

$$\begin{aligned} \ell_{Rp}(\boldsymbol{\beta}, \mathbf{D}) = & -\frac{1}{2} \left\{ (N_T - m) \log \sum_{i=1}^N (\mathbf{y}_i - \mathbf{X}_i \boldsymbol{\beta})^T \mathbf{V}_i^{-1} (\mathbf{y}_i - \mathbf{X}_i \boldsymbol{\beta}) \right. \\ & \left. + \log \left| \sum_{i=1}^N \mathbf{X}_i^T \mathbf{V}_i^{-1} \mathbf{X}_i \right| + \sum_{i=1}^N [\log |\mathbf{V}_i|] \right\} \end{aligned}$$

REML summary

- ▶ likelihood function based on the transformed data.
- ▶ likelihood function independent of fixed effect coefficients.
- ▶ unbiased estimators for variance components.
- ▶ R functions, including **lme** and **lmer**, support both REML and ML.

Balanced Random-Coefficient Model

- ▶ Assumption 1: $\mathbf{Z} = \mathbf{X}_i = \mathbf{Z}_i$ for $i = 1, \dots, N$. (so $\mathbf{V}_i = \mathbf{V}$)
- ▶ Assumption 2: $n_i = n$ for $i = 1, \dots, N$.
- ▶ Interpretation:

$$\mathbf{y}_i = \mathbf{Z}(\boldsymbol{\beta} + \mathbf{b}_i) + \boldsymbol{\epsilon}_i,$$

where $\boldsymbol{\beta}$ is the expectation of coefficients and \mathbf{b}_i is the random part.

- ▶ Example: repeated measure of patients.

Balanced Random-Coefficient Model

- ▶ Fixed-effect coefficients:

$$\hat{\boldsymbol{\beta}}_{GLS} = \hat{\boldsymbol{\beta}}_{OLS} = (\mathbf{Z}^T \mathbf{Z})^{-1} \mathbf{Z}^T \bar{\mathbf{y}}$$

where $\mathbf{y} = N^{-1} \sum_i \mathbf{y}_i$.

- ▶ Variance:

$$\hat{\sigma}_{ML}^2 = \hat{\sigma}_{RML}^2 = \frac{1}{N(n-m)} \sum_{i=1}^N \mathbf{y}_i^T (\mathbf{I} - \mathbf{Z}(\mathbf{Z}^T \mathbf{Z})^{-1} \mathbf{Z}^T) \mathbf{y}_i$$

$$\hat{\mathbf{D}}_{ML} = \frac{1}{N \hat{\sigma}_{ML}^2} (\mathbf{Z}^T \mathbf{Z})^{-1} \mathbf{Z}^T \hat{\mathbf{E}} \hat{\mathbf{E}}^T \mathbf{Z} (\mathbf{Z}^T \mathbf{Z})^{-1} - (\mathbf{Z}^T \mathbf{Z})^{-1}$$

$$\hat{\mathbf{D}}_{RML} = \frac{1}{(N-1) \hat{\sigma}_{ML}^2} (\mathbf{Z}^T \mathbf{Z})^{-1} \mathbf{Z}^T \hat{\mathbf{E}} \hat{\mathbf{E}}^T \mathbf{Z} (\mathbf{Z}^T \mathbf{Z})^{-1} - (\mathbf{Z}^T \mathbf{Z})^{-1}$$

where $\hat{\mathbf{E}} \hat{\mathbf{E}}^T = \sum_i (\mathbf{y}_i - \mathbf{Z} \hat{\boldsymbol{\beta}})(\mathbf{y}_i - \mathbf{Z} \hat{\boldsymbol{\beta}})^T$.

Balanced Random-Coefficient Model

► Log-likelihood

$$\ell = -\frac{N}{2} \left\{ n \log \sigma^2 + \log |\mathbf{I} + \mathbf{Z}^T \mathbf{D} \mathbf{Z}| + \frac{1}{N \sigma^2} \sum_i (\mathbf{y}_i - \mathbf{Z} \boldsymbol{\beta})^T (\mathbf{I} + \mathbf{Z}^T \mathbf{D} \mathbf{Z})^{-1} (\mathbf{y}_i - \mathbf{Z} \boldsymbol{\beta}) \right\}$$

► Use

$$\mathbf{V}_i^{-1} = \mathbf{I} - \mathbf{Z}(\mathbf{D}^{-1} + \mathbf{Z}^T \mathbf{Z})^{-1} \mathbf{Z}^T$$
$$\mathbf{V}_i^{-1} \mathbf{Z} = \mathbf{Z}(\mathbf{D}^{-1} + \mathbf{Z}^T \mathbf{Z})^{-1} \mathbf{D}^{-1}$$

Fitting Linear Mixed-effect Models in R

- ▶ Datasets
 - ▶ Download from the author's GitHub repository.
 - ▶ <https://github.com/eugenedemidenko/mixedmodels>
 - ▶ datasets are stored in .txt files in **Data/MixedModels** folder.
- ▶ Packages
 - ▶ `lme` function from `nlme` library.
 - ▶ `lmer` function from `lme4` library.
 - ▶ `lme` supports more covariance structures.
 - ▶ `lmer` has better scalability.

Load Dataset

Load the height-weight datasets from Family.txt file.

```
1 data = read.table("../Data/MixedModels/Chapter02/Family.txt",  
2                   header=T, stringsAsFactors=F)  
3 head(data)
```

	Height	Weight	Sex	ParentChild	Age	FamilyID
1	67.0	215	1	1	75	1
2	64.0	155	0	1	63	1
3	63.5	145	0	0	29	1
4	71.0	227	1	0	26	1
5	61.0	120	0	0	24	1
6	68.0	220	1	0	22	1

Fit LME with lme() Function

```
1 library(nlme)
2 fit.lme = lme(fixed=Weight~Height, random=~1|FamilyID, data=data)
3 fit.lme
```

- ▶ fixed argument specifies the fixed effect model. In this example, it is the linear regression of Weight against Height.
- ▶ random argument specifies the random effect model.
 - ▶ ~1 specifies that the random effect is on the intercept.
 - ▶ FamilyID specifies the group variable.

Fit LME with lme() Function

Linear mixed-effects model fit by REML

Data: data

Log-restricted-likelihood: -331.6369

Fixed: Weight ~ Height

(Intercept)	Height
-206.832149	5.345309

Random effects:

Formula: ~1 | FamilyID

(Intercept)	Residual
14.07057	24.7059

Number of Observations: 71

Number of Groups: 18

- ▶ Default estimation method: REML
- ▶ Coefficients from the fixed-effect model.
- ▶ Standard deviation for the random-effect coefficients.

Fit LME with lme() Function

To use ML, need to specify the method argument:

```
1 fit.lme = lme(fixed=Weight~Height, random=~1|FamilyID,  
2             method="ML", data=data)  
3 fit.lme
```

Fit LME with lme() Function

Linear mixed-effects model fit by maximum likelihood

Data: data

Log-likelihood: -334.7041

Fixed: Weight ~ Height

(Intercept) Height

-205.015367 5.319309

- ▶ (a) Different likelihood values.
- ▶ (b) Different fixed-effect coefficients.
- ▶ (c) Different variance parameters.

Random effects:

Formula: ~1 | FamilyID

(Intercept) Residual

StdDev: 13.34261 24.50155

Number of Observations: 71

Number of Groups: 18

Why?

Fit LME with lme() Function

Linear mixed-effects model fit by maximum likelihood

```
Data: data
Log-likelihood: -334.7041
Fixed: Weight ~ Height
(Intercept)      Height
-205.015367      5.319309
```

Random effects:

```
Formula: ~1 | FamilyID
          (Intercept) Residual
StdDev:    13.34261  24.50155
```

Number of Observations: 71

Number of Groups: 18

- ▶ (a) Different likelihood values.
- ▶ (b) Different fixed-effect coefficients.
- ▶ (c) Different variance parameters.

Why?

- ▶ (a) and (c): the use of restricted likelihood.
- ▶ (b): same formula for $\hat{\beta}_{GLS}$ but with different \hat{D} 's.

Fit LME with `lme()` Function

To get more output information, we can call

```
1 summary(fit.lme)
```

```
Linear mixed-effects model fit by maximum likelihood
```

```
Data: data
```

```
      AIC      BIC    logLik  
677.4082 686.4589 -334.7041
```

```
Random effects:
```

```
Formula: ~1 | FamilyID
```

```
(Intercept) Residual
```

```
StdDev:    13.34261 24.50155
```

```
Fixed effects: Weight ~ Height
```

	Value	Std.Error	DF	t-value	p-value
(Intercept)	-205.01537	54.08819	52	-3.790390	4e-04
Height	5.31931	0.78212	52	6.801126	0e+00

```
Correlation:
```

```
(Intr)
```

```
Height -0.997
```

```
Standardized Within-Group Residuals:
```

	Min	Q1	Med	Q3	Max
	-2.10329358	-0.54475601	-0.08698002	0.36755416	3.64634517

```
Number of Observations: 71
```

```
Number of Groups: 18
```

Fit LME with `lme()` Function

- ▶ height coefficient is random

Fit LME with `lme()` Function

- ▶ height coefficient is random

```
1 lme(fixed=Weight ~ Height, random=~1+Height | FamilyID, data=data)
```

- ▶ group indicated by FamilyID and Sex

Fit LME with lme() Function

- ▶ height coefficient is random

```
1 lme(fixed=Weight ~ Height, random=~1+Height|FamilyID, data=data)
```

- ▶ group indicated by FamilyID and Sex

```
1 lme(fixed=Weight ~ Height, random=~1|FamilyID/Sex, data=data)
```

Fit LME with lmer() Function

```
1 library(lme4)
2 fit.lme = lmer(Weight~Height+(1|FamilyID), data=data)
3 fit.lme
```

Linear mixed model fit by REML ['lmerMod']

Formula: Weight ~ Height + (1 | FamilyID)

Data: data

REML criterion at convergence: 663.2737

Random effects:

Groups	Name	Std.Dev.
--------	------	----------

	FamilyID (Intercept)	14.07
--	----------------------	-------

	Residual	24.71
--	----------	-------

Number of obs: 71, groups: FamilyID, 18

Fixed Effects:

(Intercept)	Height
-------------	--------

-206.832	5.345
----------	-------

Fit LME with lmer() Function

Do not forget "()" for random effects

```
1 fit.lme.wrong = lmer(Weight ~ Height + 1 | FamilyID, data=data)
2 fit.lme.wrong
```

Linear mixed model fit by REML ['lmerMod']

Formula: Weight ~ Height + 1 | FamilyID

Data: data

REML criterion at convergence: 697.5127

Random effects:

Groups	Name	Std.Dev.	Corr
FamilyID	(Intercept)	33.4422	
	Height	0.5802	-1.00
Residual		33.5836	

Number of obs: 71, groups: FamilyID, 18

Fixed Effects:

(Intercept)

160.9

optimizer (nloptwrap) convergence code: 0 (OK) ; 0 optimizer warnings; 2 lme4 warnings

Fit LME with `lmer()` Function

Use REML argument to choose the estimation method.

```
1 fit.lme = lmer(Weight ~ Height + (1 | FamilyID), REML=F, data=data)
2 fit.lme
```

Linear mixed model fit by maximum likelihood [`'lmerMod'`]

Formula: `Weight ~ Height + (1 | FamilyID)`

Data: `data`

AIC	BIC	logLik	deviance	df.resid
677.4082	686.4589	-334.7041	669.4082	67

Random effects:

Groups	Name	Std.Dev.
	FamilyID (Intercept)	13.34
	Residual	24.50

Number of obs: 71, groups: FamilyID, 18

Fixed Effects:

(Intercept)	Height
-205.015	5.319

Fit LME with lmer() Function

- ▶ height coefficient is random

```
1 lmer(Weight ~ Height + (1 + Height | FamilyID), data = data)
```

- ▶ Two-way group structure

```
1 lmer(Weight ~ Height + (1 | FamilyID / Sex), data = data)
```

- ▶ Force independent random effects

```
1 lmer(Weight ~ Height + (1 + Height || FamilyID), data = data)
```


Maximization Algorithms

To maximize $\ell(\boldsymbol{\theta})$, the iterative algorithms update the value of $\boldsymbol{\beta}$ by

$$\boldsymbol{\theta}^{(t+1)} = \boldsymbol{\theta}^{(t)} + \lambda_t \mathbf{H}_t^{-1} \nabla_{\boldsymbol{\theta}} \ell.$$

- ▶ Newton Raphson: $\mathbf{H} = -\nabla \ell \nabla^T$ is the negative Hessian matrix of ℓ .
- ▶ Fisher scoring: $\mathbf{H} = -\mathbb{E}[\nabla \ell \nabla]$ is the negative information matrix.
- ▶ EM algorithm: \mathbf{H} is some positive definite matrix.

A Note on EM Algorithm

EM algorithm is used to maximize the marginal likelihood function when missing data exists.

$$\max_{\boldsymbol{\theta}} \ell(X | \boldsymbol{\theta}) = \log \int L(X, Y | \boldsymbol{\theta}) dY$$

- ▶ Expectation step: compute the expected complete log-likelihood function given the observed data X and the parameter.

$$Q(\boldsymbol{\theta} | \boldsymbol{\theta}^{(i)}) = \int \ell(X, Y | \boldsymbol{\theta}) p(Y | X, \boldsymbol{\theta}^{(i)}) dY$$

- ▶ Maximization step: maximize the expected log-likelihood function.

$$\boldsymbol{\theta}^{(i+1)} = \arg \max_{\boldsymbol{\theta}} Q(\boldsymbol{\theta} | \boldsymbol{\theta}^{(i)})$$