

# Overview

## 1 Introduction

- Longitudinal Data
- Variation and Correlation
- Different Approaches

## 2 Mixed Models

- Linear Mixed Models
- Generalized Linear Mixed Models

## 3 Marginal Models

## 4 Transition Models

## 5 Further Topics

# The General Linear Mixed Model (LMM)

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{b} + \boldsymbol{\varepsilon}, \quad (1)$$

- $\mathbf{X}$  and  $\mathbf{Z}$  known design matrices ( $n \times p$  and  $n \times r$ )
- $\boldsymbol{\beta}$  vector of unknown fixed parameters
- $\mathbf{b}$  vector of random effects
- $\boldsymbol{\varepsilon}$  vector of unobservable random errors

**Assumptions:** Independence of  $\mathbf{b}$  and  $\boldsymbol{\varepsilon}$ , and

$$E \begin{pmatrix} \mathbf{b} \\ \boldsymbol{\varepsilon} \end{pmatrix} = \begin{pmatrix} \mathbf{0} \\ \mathbf{0} \end{pmatrix} \quad \text{and} \quad \text{Cov} \begin{pmatrix} \mathbf{b} \\ \boldsymbol{\varepsilon} \end{pmatrix} = \begin{pmatrix} \mathbf{D} & \mathbf{0} \\ \mathbf{0} & \mathbf{R} \end{pmatrix} \quad (2)$$

Typically, also multivariate normality of  $(\mathbf{b}', \boldsymbol{\varepsilon}')'$  is assumed. (3)

# The Longitudinal Linear Mixed Model

The longitudinal linear mixed model is typically given on the level of the subject.

$$\mathbf{Y}_i = \mathbf{X}_i\boldsymbol{\beta} + \mathbf{Z}_i\mathbf{b}_i + \boldsymbol{\varepsilon}_i, \quad \mathbf{b}_i \sim N(\mathbf{0}, \mathbf{D}_0), \quad \boldsymbol{\varepsilon}_i \sim N(\mathbf{0}, \mathbf{R}_i) \quad (4)$$

where the subscript  $i$  indicates subject-specific quantities ([Laird & Ware, 1982](#)).

A typical longitudinal model is the growth-curve model

$$Y_{ij} = \beta_0 + \beta_1 t_{ij} + b_{0i} + b_{1i} t_{ij} + \varepsilon_{ij}, \quad j = 1, \dots, n_i, \quad i = 1, \dots, I,$$

with

$$\begin{pmatrix} b_{0i} \\ b_{1i} \end{pmatrix} \stackrel{iid}{\sim} N\left(\mathbf{0}, \begin{pmatrix} \sigma_1^2 & \sigma_{12} \\ \sigma_{12} & \sigma_2^2 \end{pmatrix}\right) \quad \text{independent of} \quad \varepsilon_{ij} \stackrel{iid}{\sim} N(0, \sigma^2).$$

Here,  $\mathbf{D}_0 = \begin{pmatrix} \sigma_1^2 & \sigma_{12} \\ \sigma_{12} & \sigma_2^2 \end{pmatrix}$  and  $\mathbf{R}_i = \sigma^2 \mathbf{I}_{n_i}$ .

Example: A two-stage analysis → blackboard

# The Longitudinal Linear Mixed Model

- Longitudinal linear mixed model (LLMM) special case of LMM (1).  
Let  $\mathbf{D}$ ,  $\mathbf{R}$  and  $\mathbf{Z}$  be block-diagonal with blocks  $\mathbf{D}_0$ ,  $\mathbf{R}_i$  and  $\mathbf{Z}_i$ , respectively.  
Stack  $\mathbf{Y}_i$ ,  $\mathbf{X}_i$ ,  $\mathbf{b}_i$  and  $\varepsilon_i$  to obtain  $\mathbf{Y}$ ,  $\mathbf{X}$ ,  $\mathbf{b}$  and  $\varepsilon$ .
- LLMM assumes that  $\mathbf{Y}$  can be divided into independent subvectors  $\mathbf{Y}_i$  for the  $i$ th subject.
- LMM (1) more general - allows for inclusion of other model components such as **smooth functions** modeled using **penalized splines** in the mixed model framework. → Loss of this independence, but very flexible models.
- LMM (1) can also incorporate additional clusters (e.g. growth curves for children clustered in schools), or nested structures of random effects.

# Excursus: Penalized Splines in a Nutshell

- Consider the **nonparametric regression problem**

$$Y_i = m(x_i) + \varepsilon_i, \quad i = 1, \dots, n,$$

with **unknown smooth function**  $m(\cdot)$ .

- Approximate  $m(\cdot)$  by a linear combination of (many) **spline basis functions**, e.g. truncated polynomials,

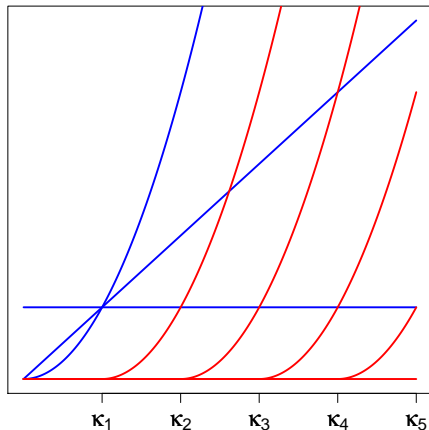
$$m(x) \approx \sum_{j=0}^q \beta_j x^j + \sum_{j=1}^K b_j (x - \kappa_j)_+^q,$$

where  $\kappa_1, \dots, \kappa_K$  is a sequence of knots, and  $u_+^q = (\max\{0, u\})^q$ .

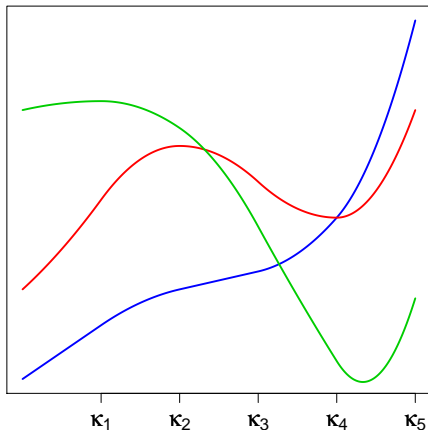
- This approximation gives a piecewise polynomial function of degree  $q$  with certain smoothness properties (continuity of  $(q - 1)$ th derivative).

# Excursus: Penalized Splines in a Nutshell

Quadratic Spline Basis



Three Example Functions



# Excursus: Penalized Splines in a Nutshell

- Estimation using a **regularization penalty**, to avoid overly wiggly function:

$$(\mathbf{Y} - \mathbf{X}\boldsymbol{\beta} - \mathbf{Z}\mathbf{b})'(\mathbf{Y} - \mathbf{X}\boldsymbol{\beta} - \mathbf{Z}\mathbf{b}) + \frac{1}{\lambda} \mathbf{b}'\mathbf{b} \rightarrow \min_{\boldsymbol{\beta}, \mathbf{b}},$$

- $\mathbf{X}$  and  $\mathbf{Z}$  design matrices for the spline basis functions  $x^j$  resp.  $(x - \kappa_j)_+^q$
  - $\lambda$  a smoothing parameter.
  - Penalized least squares estimation equivalent to estimation in the LMM with
    - **fixed effects**  $\beta_0, \dots, \beta_q$ : subspace of polynomial functions (degree  $q$ )
    - **random effects**  $b_1, \dots, b_K \stackrel{iid}{\sim} N(0, \sigma_b^2)$ : deviations from subspace.
- Can estimate **smoothing parameter**  $\lambda = \sigma_b^2 / \sigma_\varepsilon^2$  **data driven**.
- Analogous for other basis choices (e.g. B-Splines) and for spatial effects, interaction surfaces, varying coefficients, ...  $\rightarrow$  can combine with random effects for longitudinal data in a general LMM (**additive mixed model**).
  - Examples: smooth dose-response functions for covariates, smooth profiles over time, time-varying effects, spatial surfaces (spatial information on subjects), .... More detailed information: **Ruppert, Wand & Carroll, 2003**.

# $D$ , $R$ and $V$

- $D$ ,  $R$  and  $Z$  together imply the correlation structure for  $\mathbf{Y}$ , as  $\mathbf{V} = \text{Cov}(\mathbf{Y}) = \mathbf{ZDZ}' + \mathbf{R}$ .
- $\mathbf{Zb}$  with  $D$  models differences between subjects, e.g. in their average level and average evolution in time.
- $R$  models the residual serial correlation not explained by  $\mathbf{Zb}$ , as well as measurement error.



# Some common models and implications for $\mathbf{V}$

The **random intercept** model

$$Y_{ij} = \mathbf{x}_{ij}'\boldsymbol{\beta} + b_i + \varepsilon_{ij}, \quad b_i \stackrel{iid}{\sim} N(0, \sigma_b^2), \quad \varepsilon_{ij} \stackrel{iid}{\sim} N(0, \sigma_\varepsilon^2)$$

with  $\mathbf{R} = \sigma_\varepsilon^2 \mathbf{I}_n$ ,  $\mathbf{D} = \sigma_b^2 \mathbf{I}_I$ ,  $\mathbf{D}_0 = \sigma_b^2$  and  $\mathbf{Z}_i = (1, \dots, 1)'$  implies  $\mathbf{V}$  is block-diagonal, with the  $i$ th block for the  $i$ th subject equal to

$$\begin{pmatrix} \sigma_\varepsilon^2 + \sigma_b^2 & \sigma_b^2 & \dots & \sigma_b^2 \\ \sigma_b^2 & \sigma_\varepsilon^2 + \sigma_b^2 & \dots & \sigma_b^2 \\ \vdots & \ddots & \ddots & \vdots \\ \sigma_b^2 & \sigma_b^2 & \dots & \sigma_\varepsilon^2 + \sigma_b^2 \end{pmatrix}.$$

This is called **compound symmetry**. In a marginal model, this is often written with  $\sigma^2$  for  $\sigma_\varepsilon^2 + \sigma_b^2$  and  $\rho\sigma^2$  for  $\sigma_b^2$ . (Then, the correlation  $\rho$  could be negative.)

This implies that **all observations** on a subject are equally correlated, i.e. it assumes that observations on the same subject closer in time are not more similar than observations further away in time. (E.g. blood marker with short half-life.)

# Some common models and implications for $V$

The **growth-curve** model

$$Y_{ij} = \beta_0 + \beta_1 t_{ij} + b_{0i} + b_{1i} t_{ij} + \varepsilon_{ij}, \quad b_i \stackrel{iid}{\sim} N\left(\mathbf{0}, \begin{pmatrix} \sigma_1^2 & \sigma_{12} \\ \sigma_{12} & \sigma_2^2 \end{pmatrix}\right), \quad \varepsilon_{ij} \stackrel{iid}{\sim} N(0, \sigma_\varepsilon^2)$$

with  $\mathbf{D}$  block-diagonal with blocks  $\begin{pmatrix} \sigma_1^2 & \sigma_{12} \\ \sigma_{12} & \sigma_2^2 \end{pmatrix}$ ,  $\mathbf{Z}_i = \begin{pmatrix} 1 & \dots & 1 \\ t_{i1} & \dots & t_{i n_i} \end{pmatrix}'$  and  $\mathbf{R} = \sigma_\varepsilon^2 \mathbf{I}_n$ , implies

$$\begin{aligned} \text{Var}(Y_{ij}) &= \sigma_1^2 + 2\sigma_{12}t_{ij} + \sigma_2^2 t_{ij}^2 + \sigma_\varepsilon^2 \quad \text{and} \\ \text{Cov}(Y_{ij}, Y_{ik}) &= \sigma_1^2 + \sigma_{12}t_{ij} + \sigma_{12}t_{ik} + \sigma_2^2 t_{ij}t_{ik}, \quad j \neq k. \end{aligned}$$

In particular, the variance function is **quadratic** in time. Adding a quadratic term  $b_{2i}t_{ij}^2$  would make it a forth-order polynomial in time. (Extrapolation??)

# Residual Covariance Structure

Often, the **conditional independence model**  $\mathbf{R} = \sigma_\varepsilon^2 \mathbf{I}_n$  is chosen, which implies that the responses of a subject  $i$  are independent, given  $\beta$  and  $\mathbf{b}_i$ . This may be unrealistically simplistic. If the covariance structure can not be meaningfully explained by random effects alone, more general forms of  $\mathbf{R}$  should be chosen.

$\mathbf{R}$  is typically chosen to be block-diagonal with blocks depending on  $i$  only through  $n_i$  and the time points  $t_{ij}, j = 1, \dots, n_i$ . Usually, it is modeled as

$$\text{Cov}(\varepsilon_{ij}, \varepsilon_{ik}) = \sigma^2 + \tau^2 g(|t_{ij} - t_{ik}|)$$

with a suitable decreasing function  $g(\cdot)$  with  $g(0) = 1$  and  $\lim_{u \rightarrow \infty} g(u) = 0$ .

Two frequently used functions are the **exponential** and **Gaussian correlation functions**  $g(u) = \exp(-u)$  and  $g(u) = \exp(-\phi u^2)$  (with  $\phi > 0$ ).

In applications, however, the effect of serial correlation is often dominated by the combination of random effects and measurement error (also causing estimation problems when all are present in the model).

# Conditional and Marginal Perspectives

Conditional perspective on the mixed model:

$$\mathbf{Y}|\mathbf{b} \sim N(\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{b}, \sigma^2 \mathbf{I}_n), \quad \mathbf{b} \sim N(\mathbf{0}, \mathbf{D}). \quad (5)$$

Interpretation: Random effects are subject-specific effects on the mean that vary within the population and are estimated subject to a regularization constraint.

Marginal perspective on the mixed model:

$$\mathbf{Y} \sim N(\mathbf{X}\boldsymbol{\beta}, \mathbf{V}), \quad (6)$$

where  $\mathbf{V} = \text{Cov}(\mathbf{Y}) = \mathbf{Z}\mathbf{D}\mathbf{Z}' + \mathbf{R}$ .

Interpretation: The random effects induce a correlation structure and thus enable a proper statistical analysis of the correlated data.

# Inference in the Linear Mixed Model

What is the likelihood for the linear mixed model (1)-(3)?

Inference is usually based on the marginal model (6) for  $\mathbf{Y}$  (unless Bayesian methods are used).

Then, the unknown parameters are  $(\beta, \theta)$ , where  $\theta$  contains the unknown parameters in  $\mathbf{D}$  and  $\mathbf{R}$  (and thus  $\mathbf{V}$ ).

(5) implies (6), but they are not equivalent!

**Example:** random intercept + independent errors  $\Rightarrow$  compound symmetry for  $\mathbf{V}$

$\sigma_b^2$  covariance in (6)  $\rightarrow$  could be negative

$\sigma_b^2$  variance in (5)  $\rightarrow$  has to be non-negative

And (6) does not include  $\mathbf{b}$ . Not as much of a problem if interest is in  $(\beta, \theta)$ , but does not work if interest is in  $\mathbf{b}$ .

## Estimation of $\beta$

For a given  $\theta$ , (6) is simply a general linear model. We can estimate  $\beta$  using generalized least squares

$$(\mathbf{Y} - \mathbf{X}\beta)' \mathbf{V}^{-1} (\mathbf{Y} - \mathbf{X}\beta) \rightarrow \min_{\beta}$$

which gives

$$\hat{\beta} = (\mathbf{X}' \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}' \mathbf{V}^{-1} \mathbf{Y}. \quad (7)$$

(We assume that the inverses  $\mathbf{V}^{-1}$  and  $(\mathbf{X}' \mathbf{V}^{-1} \mathbf{X})^{-1}$  exist. Generalizations use generalized inverses.)

This is also the maximum likelihood (ML) estimator when we have the normality assumption (3): The log-likelihood for  $\beta$  and  $\theta$  from the marginal model (6) is

$$\ell(\beta, \theta) = -\frac{n}{2} \log(2\pi) - \frac{1}{2} \log |\mathbf{V}| - \frac{1}{2} (\mathbf{Y} - \mathbf{X}\beta)' \mathbf{V}^{-1} (\mathbf{Y} - \mathbf{X}\beta). \quad (8)$$

$$\frac{\partial}{\partial \beta} \ell(\beta, \theta) = \mathbf{X}' \mathbf{V}^{-1} (\mathbf{Y} - \mathbf{X}\beta) \stackrel{!}{=} \mathbf{0} \quad \Rightarrow \quad \hat{\beta} = (\mathbf{X}' \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}' \mathbf{V}^{-1} \mathbf{Y}.$$

# Estimation of $\beta$

$\hat{\beta}$  is also the BLUE, the best linear unbiased estimator ([Gauss-Markov theorem](#)):

Let  $\tilde{\beta} = \mathbf{a} + \mathbf{B}\mathbf{y}$  be another linear unbiased estimator of  $\beta$ .

$$\Rightarrow E(\tilde{\beta}) = \mathbf{a} + \mathbf{B}\mathbf{X}\beta = \beta \Rightarrow \mathbf{a} = \mathbf{0} \text{ and } \mathbf{B}\mathbf{X} = \mathbf{I}_p.$$

$\Rightarrow$  For all  $\mathbf{c} \in \mathbb{R}^p$

$$\begin{aligned} & \text{Var}(\mathbf{c}'\tilde{\beta}) - \text{Var}(\mathbf{c}'\hat{\beta}) \\ &= \mathbf{c}'\mathbf{B}\mathbf{V}\mathbf{B}'\mathbf{c} - \mathbf{c}'(\mathbf{X}'\mathbf{V}^{-1}\mathbf{X})^{-1}\mathbf{X}'\mathbf{V}\mathbf{X}(\mathbf{X}'\mathbf{V}^{-1}\mathbf{X})^{-1}\mathbf{c} \\ &= \mathbf{c}'\mathbf{B}\mathbf{V}\mathbf{B}'\mathbf{c} - \mathbf{c}'\mathbf{B}\mathbf{X}(\mathbf{X}'\mathbf{V}^{-1}\mathbf{X})^{-1}\mathbf{X}'\mathbf{B}'\mathbf{c} \\ &= \mathbf{c}'\mathbf{B}\mathbf{V}^{1/2}[\mathbf{I}_n - \mathbf{V}^{-1/2}\mathbf{X}(\mathbf{X}'\mathbf{V}^{-1}\mathbf{X})^{-1}\mathbf{X}'\mathbf{V}^{-1/2}]\mathbf{V}^{1/2}\mathbf{B}'\mathbf{c} \geq 0 \end{aligned}$$

as the middle matrix is a projection matrix and thus positive semi-definite.

As  $\hat{\beta}$  is unbiased with smallest variance, it is the BLUE.

With the normality assumption (3),  $\hat{\beta}$  can even be shown to be the best unbiased estimator.

# Estimation of $\theta$

## Maximum likelihood estimation

Substitute  $\hat{\beta} = \hat{\beta}(\theta)$  into (8)  $\rightarrow$  profile log-likelihood  $\ell(\hat{\beta}(\theta), \theta)$ .

Maximize to obtain ML estimate for  $\theta$ .

No closed form solution in general  $\rightarrow$  numerical maximization.

## Restricted maximum likelihood estimation

Maximum likelihood estimates for variances are biased downward.

**Example:**  $Y_i \stackrel{iid}{\sim} N(\mu, \sigma^2)$ , ML estimator  $\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (Y_i - \bar{Y})^2$ , but  $E(\hat{\sigma}^2) = \frac{n-1}{n} \sigma^2$ .

$\rightarrow$  Everyone uses  $\hat{\sigma}^2 = \frac{1}{n-1} \sum_{i=1}^n (Y_i - \bar{Y})^2$ .

This is the restricted maximum likelihood estimate.

It accounts for estimation of  $\mu$  by  $\bar{Y}$ .



# Restricted Maximum Likelihood Estimation of $\theta$

**Idea:** Use likelihood not of  $\mathbf{Y}$ , but of  $(n - p)$  linearly independent error contrasts  $\mathbf{AY}$  with distribution independent of unknown  $\beta$ . (Patterson & Thompson, 1971).

Split the log-likelihood into two parts for  $\mathbf{AY}$  (for  $\theta$ ) and  $\mathbf{BY}$  (for  $\beta$ ) with

- ①  $\mathbf{A}$  of rank  $n - p$  and  $\mathbf{B}$  of rank  $p$ .
- ② The two parts are statistically independent, i.e. (here:)  $\text{Cov}(\mathbf{AY}, \mathbf{BY}) = \mathbf{0}$ .
- ③  $\mathbf{AY}$  are error contrasts, i.e.  $E(\mathbf{AY}) = \mathbf{0}$ .
- ④  $\mathbf{BX}$  is of rank  $p$  ( $\Rightarrow \mathbf{BY}$  estimates one-to-one function of  $\beta$ ).

Suitable matrices:  $\mathbf{A} = \mathbf{I}_n - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'$  and  $\mathbf{B} = (\mathbf{X}'\mathbf{V}^{-1}\mathbf{X})^{-1}\mathbf{X}'\mathbf{V}^{-1}$ .

→ Alternative name **residual maximum likelihood estimation**.

From the first two points, the log-likelihood is  $\ell = \ell' + \ell^*$  (up to an additive constant),  $\ell'$  for  $\mathbf{AY}$  and  $\ell^*$  for  $\mathbf{BY}$ .

# Restricted Maximum Likelihood Estimation of $\theta$

Maximizing  $\ell^*$  (from  $\mathbf{B}\mathbf{Y}$ ) with respect to  $\beta$  gives again the MLE/BLE  $\hat{\beta}$  (depending on  $\theta$ ).

Maximizing  $\ell'$  (from  $\mathbf{A}\mathbf{Y}$ ) with respect to  $\theta$ :

Harville (1974) showed that the restricted log-likelihood  $\ell'$  is (up to an additive constant) independent of the precise error contrasts used.

He derived the restricted log-likelihood, which is (up to a constant)

$$\begin{aligned} l(\theta) &= -\frac{1}{2} \log |\mathbf{V}| - \frac{1}{2} \log |\mathbf{X}\mathbf{V}^{-1}\mathbf{X}| - \frac{1}{2} (\mathbf{Y} - \mathbf{X}\hat{\beta})' \mathbf{V}^{-1} (\mathbf{Y} - \mathbf{X}\hat{\beta}) \\ &= \ell(\hat{\beta}(\theta), \theta) - \frac{1}{2} \log |\mathbf{X}\mathbf{V}^{-1}\mathbf{X}| \end{aligned}$$

(Proof: blackboard)

Maximize  $\rightarrow$  restricted maximum likelihood (REML) estimate for  $\theta$ .

**Example:** In a linear model, this gives the usual estimate for  $\sigma^2$  (dividing by  $n - p$ ). It takes into account the loss in degrees of freedom from estimation of  $\beta$ .

# Restricted Maximum Likelihood Estimation of $\theta$

The restricted likelihood is also the conditional likelihood for  $\theta$  given the sufficient statistic  $\mathbf{BY}$  for  $\beta$  and it can be argued that this entails no information loss (Smyth & Verbyla, 1996).

The concept that  $\mathbf{BY}$  contains no information about  $\theta$  absent knowledge about  $\beta$  has also been called marginal sufficiency (Sprott, 1975; Harville, 1977).

# Implementation of ML or REML

Lindstrom & Bates, JASA, 1988 derive the first- and second-order derivatives of log-likelihood and restricted log-likelihood. They give details of implementation of a Newton-Raphson Algorithm for estimation.

Note that the requirements for  $\theta$  in the marginal model ( $\mathbf{V}$  positive (semi-)definite) are different from those in the conditional model ( $\mathbf{D}_0$  and  $\mathbf{R}_i$  positive (semi-)definite). Thus, software packages maximize for  $\theta$  often over a larger parameter space than the conditional model implies. E.g. SAS proc mixed (at least v6.12) by default only requires the diagonal elements of  $\mathbf{D}_0$  and  $\mathbf{R}_i$  to be positive. Also, be aware of other numerical issues. For example, R lme() maximizes the (restricted) log-likelihood with respect to the scaled logarithm of the variances, and thus can never find a maximum at zero (Pinheiro & Bates, 2000, section 2.2).

An alternative is the EM algorithm, treating the random effects as missing data, which, however, can be slow to converge. R lme() uses a hybrid of EM and then Newton-Raphson.

## EBLUE and EBLUP

After estimation of  $\boldsymbol{\theta}$ , we can substitute  $\hat{\boldsymbol{\theta}}$  into the formula  $\hat{\boldsymbol{\beta}} = \hat{\boldsymbol{\beta}}(\boldsymbol{\theta})$  to obtain an estimate. This is the **estimated BLUE (EBLUE)**.

Often it is also of interest to obtain predictions of  $\mathbf{b}$  (e.g. find “unusual” subjects). Then, we have to go back to the conditional model formulation (5). The best linear unbiased predictor (**BLUP**) for  $\mathbf{b}$  is (see next slide)

$$\hat{\mathbf{b}} = \mathbf{DZ}'\mathbf{V}^{-1}(\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}).$$

The combined BLUE and BLUP (or collectively referred to as BLUPs) can also be written as

$$\begin{pmatrix} \hat{\boldsymbol{\beta}} \\ \hat{\mathbf{b}} \end{pmatrix} = (\mathbf{C}'\mathbf{R}^{-1}\mathbf{C} + \mathbf{G})^{-1}\mathbf{C}'\mathbf{R}^{-1}\mathbf{Y},$$

where  $\mathbf{C} = (\mathbf{X}|\mathbf{Z})$  and  $\mathbf{G} = \begin{pmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{D}^{-1} \end{pmatrix}$ .

The **estimated BLUP (EBLUP)** is obtained by substituting  $\hat{\boldsymbol{\theta}}$  for  $\boldsymbol{\theta}$  in  $\mathbf{D}$ ,  $\mathbf{V}$  and  $\hat{\boldsymbol{\beta}}$ .

# BLUE and BLUP

BLUE and BLUP (or collectively referred to as BLUP)  $\hat{\beta}$  and  $\hat{\mathbf{b}}$  are the solution to the simultaneous Henderson's mixed model equations ([Henderson, 1950](#))

$$\begin{aligned} \mathbf{X}'\mathbf{R}^{-1}\mathbf{X}\hat{\beta} + \mathbf{X}'\mathbf{R}^{-1}\mathbf{Z}\hat{\mathbf{b}} &= \mathbf{X}'\mathbf{R}^{-1}\mathbf{Y} \\ \mathbf{Z}'\mathbf{R}^{-1}\mathbf{X}\hat{\beta} + (\mathbf{Z}'\mathbf{R}^{-1}\mathbf{Z} + \mathbf{D}^{-1})\hat{\mathbf{b}} &= \mathbf{Z}'\mathbf{R}^{-1}\mathbf{Y} \end{aligned} \quad (9)$$

There are (at least) four different justifications for these equations ([Robinson, 1991](#))

- a) They arise when maximizing the joint density of  $\mathbf{Y}$  and  $\mathbf{b}$ .
- b) In a Bayesian approach, regard  $\beta$  as a parameter with a uniform improper prior distribution, independent of  $\mathbf{b}$ . (9) gives the posterior mode for  $\beta$  and  $\mathbf{b}$ . Thus, the resulting estimates are also called [empirical Bayes estimates](#).
- c)  $\hat{\beta}$  and  $\hat{\mathbf{b}}$  have minimum mean squared error within the class of linear unbiased estimates for  $\beta$  and  $\mathbf{b}$ . (Unbiasedness for  $\hat{\mathbf{b}}$ :  $E(\hat{\mathbf{b}}) = E(\mathbf{b}) = \mathbf{0}$ .)
- d) Goldberger's derivation: For a future observation  $\mathbf{Y}_* = \mathbf{x}_*'\beta + \mathbf{z}_*'\mathbf{b} + \epsilon_*$ , the best linear unbiased predictor is  $\mathbf{x}_*'\hat{\beta} + \mathbf{z}_*'\hat{\mathbf{b}}$ .

# Shrinkage

$\hat{\mathbf{b}}$  is a **shrinkage estimator** of  $\mathbf{b}$ , i.e. its components have less spread than the generalized least squares estimators would have when treating  $\mathbf{b}$  as fixed effects.

This seems to contradict the unbiasedness property of the BLUP. However, as [Robinson, 1991](#) points out, unbiasedness here means

$$E(\hat{\mathbf{b}}) = E(\mathbf{b}) = \mathbf{0},$$

not

$$E(\hat{\mathbf{b}}|\mathbf{b}) = \mathbf{b} \text{ for all } \mathbf{b}.$$

# Shrinkage

$$\begin{aligned}
 \widehat{\mathbf{Y}}_i &= \mathbf{X}_i \widehat{\boldsymbol{\beta}} + \mathbf{Z}_i \widehat{\mathbf{b}}_i \\
 &= \mathbf{X}_i \widehat{\boldsymbol{\beta}} + \mathbf{Z}_i \mathbf{D} \mathbf{Z}_i' \mathbf{V}_i^{-1} (\mathbf{Y}_i - \mathbf{X}_i \widehat{\boldsymbol{\beta}}) \\
 &= (\mathbf{V}_i - \mathbf{Z}_i \mathbf{D} \mathbf{Z}_i') \mathbf{V}_i^{-1} \mathbf{X}_i \widehat{\boldsymbol{\beta}} + \mathbf{Z}_i \mathbf{D} \mathbf{Z}_i' \mathbf{V}_i^{-1} \mathbf{Y}_i \\
 &= \mathbf{R}_i \mathbf{V}_i^{-1} \mathbf{X}_i \widehat{\boldsymbol{\beta}} + \mathbf{Z}_i \mathbf{D} \mathbf{Z}_i' \mathbf{V}_i^{-1} \mathbf{Y}_i
 \end{aligned}$$

The BLUP for  $\mathbf{Y}_i$  thus is a **weighted average** of the population-averaged profile  $\mathbf{X}_i \widehat{\boldsymbol{\beta}}$  and the individual data  $\mathbf{Y}_i$ . The observed data is shrunk towards the population-averaged profile (“borrowing of strength”).

Note that  $\mathbf{V}_i = \mathbf{R}_i + \mathbf{Z}_i \mathbf{D} \mathbf{Z}_i'$ .

- If  $\mathbf{R}_i \mathbf{V}_i^{-1}$  is large, i.e. the residual variability ( $\mathbf{R}_i$ ) is large compared to the between-subject variability ( $\mathbf{Z}_i \mathbf{D} \mathbf{Z}_i'$ ), much weight is given to the population-averaged profile  $\mathbf{X}_i \widehat{\boldsymbol{\beta}}$ .
- If  $\mathbf{R}_i \mathbf{V}_i^{-1}$  is small, the opposite is the case.



# Inference for $\beta$

For known  $\theta$ , we have

$$\hat{\beta} \sim (\beta, (\mathbf{X}'\mathbf{V}^{-1}\mathbf{X})^{-1}).$$

Given the normality assumption (3), or asymptotically,  $\hat{\beta}$  is also normal.

In practice, for unknown  $\theta$ , the covariance has to be estimated as

$$(\mathbf{X}'\hat{\mathbf{V}}^{-1}\mathbf{X})^{-1} = (\mathbf{X}'\mathbf{V}(\hat{\theta})^{-1}\mathbf{X})^{-1}.$$

Reported standard errors and confidence intervals are usually based on this estimated covariance matrix. Note that this ignores the uncertainty in  $\theta$ .

This covariance estimate is [model-based](#).

# Robust Variance Estimation

What happens if our model for  $\text{Cov}(\mathbf{Y})$  is wrong?

For example, we could assume compound symmetry with constant variance over time, but in fact variances could increase over time, and correlations might decrease with time distance.

Suppose the true covariance is  $\text{Cov}(\mathbf{Y}) = \mathbf{V}$ , but our model specifies instead  $\mathbf{W}$  as the so-called [working covariance matrix](#). Then,

$$\hat{\beta} = (\mathbf{X}'\mathbf{W}^{-1}\mathbf{X})^{-1}\mathbf{X}'\mathbf{W}^{-1}\mathbf{Y},$$

and  $\hat{\beta}$  is still consistent. However, the covariance changes to

$$\text{Cov}(\hat{\beta}) = (\mathbf{X}'\mathbf{W}^{-1}\mathbf{X})^{-1}\mathbf{X}'\mathbf{W}^{-1}\mathbf{V}\mathbf{W}^{-1}\mathbf{X}(\mathbf{X}'\mathbf{W}^{-1}\mathbf{X})^{-1}$$

and we could use  $\hat{\beta}$  together with an estimated robust covariance matrix.

# Robust Variance Estimation

A consistent estimator (for  $I \rightarrow \infty$ ) for  $\text{Cov}(\hat{\beta})$  is (Liang & Zeger, 1986)

$$\left( \sum_{i=1}^I \mathbf{x}_i' \widehat{\mathbf{W}}_i^{-1} \mathbf{x}_i \right)^{-1} \left\{ \sum_{i=1}^I \mathbf{x}_i' \widehat{\mathbf{W}}_i^{-1} (\mathbf{y}_i - \mathbf{x}_i \hat{\beta}) (\mathbf{y}_i - \mathbf{x}_i \hat{\beta})' \widehat{\mathbf{W}}_i^{-1} \mathbf{x}_i \right\} \left( \sum_{i=1}^I \mathbf{x}_i' \widehat{\mathbf{W}}_i^{-1} \mathbf{x}_i \right)^{-1}.$$

This is the **sandwich estimator** originally due to Huber, 1967 and White, 1980, also called robust or empirical variance estimator.

This estimator is valid under more general conditions, though the model-based estimator has less variance (and no asymptotic bias) when the variance model is correctly specified, i.e. the robust is less efficient than the model-based variance estimator when we do a good job of estimating the covariance.

It has also been shown that the empirical variance estimator can be highly biased when the number of clusters (here,  $I$ ) is small.

# Hypothesis Testing for $\beta$

We can use our variance estimators to construct **Wald tests** of

$$H_0 : \mathbf{L}\beta = \mathbf{0} \text{ versus } H_A : \mathbf{L}\beta \neq \mathbf{0} \quad (10)$$

or to construct confidence intervals. For this, we use that

$$T_W = (\hat{\beta} - \beta)' \mathbf{L}' \left( \mathbf{L}(\mathbf{X}'\hat{\mathbf{V}}^{-1}\mathbf{X})^{-1}\mathbf{L}' \right)^{-1} \mathbf{L}(\hat{\beta} - \beta)$$

asymptotically follows a chi-square distribution with  $\text{rank}(\mathbf{L})$  degrees of freedom. (Or use the corresponding robust covariance estimate.)

**Note:** The results given here for testing (of fixed effects or variance components) assume the longitudinal linear mixed model (4). Asymptotics are with regard to  $I \rightarrow \infty$ . Results for more general types of linear mixed models often remain elusive (Ruppert, Wand & Carroll, 2003, section 4.8).

# Hypothesis Testing for $\beta$

The Wald test statistic is based on estimated standard errors which underestimate the true variability in  $\hat{\beta}$ , as they ignore the variability introduced by estimating  $\theta$ .

This problem is often alleviated by using approximate *t*- or *F*-statistics instead:

- For a single parameter  $\beta_j$  in  $\beta$ , the distribution of  $(\hat{\beta}_j - \beta_j)/\widehat{SE}(\hat{\beta}_j)$  is approximated by a *t*-distribution.
- For general hypotheses (10), an *F*-approximation to the distribution of

$$F = \frac{(\hat{\beta} - \beta)' \mathbf{L}' \left( \mathbf{L} (\mathbf{X}' \hat{\mathbf{V}}^{-1} \mathbf{X})^{-1} \mathbf{L}' \right)^{-1} \mathbf{L} (\hat{\beta} - \beta)}{\text{rank}(\mathbf{L})}$$

is used, with numerator degrees of freedom  $\text{rank}(\mathbf{L})$ .

# Hypothesis Testing for $\beta$

The  $t$  [degrees of freedom](#) or the denominator degrees of freedom for the  $F$ -distribution are estimated from the data. Common methods include a so-called Satterthwaite-type approximation ([Satterthwaite, Psychometrika, 1941](#)). (The general idea is to match the moments of the distributions.)

For small samples, the method by [Kenward and Roger, Biometrics, 1997](#) gives better results, as it first estimates the amount by which the asymptotic covariance matrix underestimates (in a matrix sense)  $\text{Cov}(\hat{\beta})$ , and then inflates the variance-covariance matrix by this amount before using a Satterthwaite-type approach.

Computation of Satterthwaite-type degrees of freedom are computationally expensive and can add substantially to computation time.

# Hypothesis Testing for $\beta$

We could also use a **likelihood ratio test** (LRT) based on large sample theory. When  $n$  is not that large, the LRT tends to be more reliable than the Wald test.

For an LRT, we use that the test statistic

$$T_{LRT} = 2 \sup_{H_A} \ell(\beta, \theta) - 2 \sup_{H_0} \ell(\beta, \theta)$$

under the null hypothesis converges in distribution to a  $\chi^2_{\text{rank}(L)}$  variable.

We need

- nested models
- with the same covariance structure in each model.
- And don't use with REML estimation! (Due to the different fixed effects under null and alternative, the error contrasts are different and the two likelihoods thus not comparable.)

Also, be aware that the approximation can be very poor when  $n - p$  is small.

# Hypothesis Testing for $D$

Suppose now we want to test for certain variance/covariance parameters.  
Let  $K$  be the dimension of  $D_0$  and

$$D_0 = \begin{pmatrix} d_{11} & \dots & d_{1K} \\ \vdots & \ddots & \vdots \\ d_{1K} & \dots & d_{KK} \end{pmatrix} = \left( \begin{array}{c|c} D_1 & d_{1K} \\ \hline d_{1K} & \dots & d_{KK} \end{array} \right).$$

We might be interested in testing

$$H_0 : D_0 = \left( \begin{array}{c|c} D_1 & 0 \\ \hline 0 & \dots & 0 \end{array} \right) \text{ with } D_1 \text{ positive-definite } (K-1 \times K-1) \text{ versus}$$

$$H_A : D_0 \text{ positive-semidefinite } (K \times K).$$

An example would be to test for a random intercept model versus a growth-curve model ( $K = 2$ ), or for no random effect versus a random intercept ( $K = 1$ ).



# Hypothesis Testing for $D$

This case violates one of the standard regularity assumptions, that the tested parameters be in the [interior of the parameter space](#). Instead,  $d_{KK}$  as a variance has to be non-negative and  $d_{KK} = 0$  thus is on the [boundary](#) of the parameter space. (Actually,  $D_0$  and  $D_1$  have to fulfill more complicated restrictions to be positive-(semi)definite.)

One can show that for the above hypothesis,  $T_{LRT}$  under the null converges in distribution to a variable with a 0.5:0.5 mixture between a  $\chi^2_{K-1}$  and a  $\chi^2_K$  distribution ([Self & Liang, JASA, 1987](#), [Stram & Lee, Biometrics, 1994, 1995](#), [Giampaoli & Singer, JSPI, 2009](#)).

The asymptotics assume the longitudinal linear mixed model (4) with  $I \rightarrow \infty$ . [Crainiceanu & Ruppert \(JRSS-B, 2004\)](#) show that if  $\mathbf{Y}$  cannot be subdivided into a large number of independent subvectors  $\mathbf{Y}_i$ , this mixture distribution is typically a bad approximation.

# Hypothesis Testing for $D$

Crainiceanu & Ruppert (JRSS-B, 2004) also derive the exact distribution for the case  $K = 1$  (for the general LMM (1)). This is implemented in the R-package `RLRsim`. Some extensions can be found in Greven et al (JCGS, 2008) and Scheipl, Greven & Küchenhoff (CSDA, 2008).

Note also that the  $\chi^2_{K-1} : \chi^2_K$  mixture distribution given above is not valid, if a nuisance parameter is also on the boundary (e.g. if, when testing for a random slope, the random intercept is also (near) zero). Self & Liang, JASA, 1987 show that then the whole geometry of the problem changes.

For more complex hypotheses for  $D_0$ , more complex chi-square mixtures can arise (Self & Liang, JASA, 1987, Stram & Lee, Biometrics, 1994, 1995).

When testing only for variance components, one can also use a restricted LRT (RLRT) after REML estimation. The asymptotics are the same.

# Hypothesis Testing for $D$

Molenberghs & Verbeke, *American Statistician*, 2007 look at Score-Tests and Wald-Tests in addition to LRTs and discuss how the same asymptotics apply to these two tests. However, they are more difficult to compute in the constraint case, as they require additional derivatives as well as numerical constraint optimization.

# Precision of the BLUPs

Remember

$$\begin{pmatrix} \hat{\beta} \\ \hat{\mathbf{b}} \end{pmatrix} = (\mathbf{C}'\mathbf{R}^{-1}\mathbf{C} + \mathbf{G})^{-1}\mathbf{C}'\mathbf{R}^{-1}\mathbf{Y},$$

where  $\mathbf{C} = (\mathbf{X}|\mathbf{Z})$  and  $\mathbf{G} = \begin{pmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{D}^{-1} \end{pmatrix}$ .

The common variance-covariance matrix can be shown to be (blackboard)

$$\text{Cov} \begin{pmatrix} \hat{\beta} \\ \hat{\mathbf{b}} - \mathbf{b} \end{pmatrix} = (\mathbf{C}'\mathbf{R}^{-1}\mathbf{C} + \mathbf{G})^{-1}.$$

Note that  $\text{Cov}(\hat{\mathbf{b}} - \mathbf{b})$  is used rather than  $\text{Cov}(\hat{\mathbf{b}})$ , as the latter does not recognize the random variation in  $\mathbf{b}$ .

Estimated versions can be used to construct approximate confidence/prediction intervals for expressions involving the BLUPs (such as  $\mathbf{Y}_i$  or new observations  $\mathbf{Y}_i^*$ ).

# Model Selection

Mixed model software also gives AIC and BIC values. These are based on the marginal likelihood, the likelihood from the marginal model (6). For example, the Akaike information criterion (AIC) is typically defined as

$$AIC = -2 \log \ell(\hat{\beta}, \hat{\theta}) + 2k,$$

where  $k = p + 1 + d$  is the number of unknown parameters (if  $d$  is the number of unknown parameters in  $\mathbf{D}$ ). (The AIC using the marginal restricted log-likelihood is defined analogously.)

However, be aware that the derivation of the AIC also assumes a parameter space  $\mathbb{R}^k$  and independent and identically distributed (i.i.d.) observations. The AIC thus suffers from similar problems to LRTs (boundary!) when it is used to select parameters in  $\mathbf{D}$  / random effects. It will choose the smaller model more often than you think.

# Model Selection

There are approaches to use an AIC based on the conditional log-likelihood for this case (Vaida & Blanchard, *Biometrika*, 2005; Liang, Wu & Zhou, *Biometrika*, 2008). However, they still suffer from problems and more research is needed: The first approach always chooses the larger model, and the second approach is computationally very demanding and not very robust (Greven & Kneib, *TechReport*, 2008).

# Model Diagnostics

Fully covering model diagnostics would lead to far here. A few remarks:

- **Normality assumption for the random effects:** Due to the shrinkage effect, the BLUPs  $\hat{\mathbf{b}}$  can look normal even if the true distribution for  $\mathbf{b}$  is not normal (e.g. bimodal). Verbeke & Molenberghs, 2000 give a discussion on fitting mixture distributions for  $\mathbf{b}$  instead to check for normality.
- **Residual diagnostics:** Remember that  $\text{Cov}(\mathbf{Y}_i - \mathbf{X}_i\beta) = \mathbf{V}_i$ . Thus, the residual vector  $\mathbf{r}_i = \mathbf{Y}_i - \mathbf{X}_i\hat{\beta}$  will have mean zero, but will be correlated. Diagnostics are typically based on standardized residuals  $\mathbf{r}_i^* = \mathbf{L}_i^{-1}\mathbf{r}_i$ , where  $\hat{\mathbf{V}}_i = \mathbf{L}_i\mathbf{L}_i'$  is the Cholesky decomposition. (E.g. residual plots, Mahalanobis distance between observed and predicted responses, semi-variogram.)  $\mathbf{r}_i^*$  are approximately uncorrelated with unit variance.

One could also use the subject-specific residuals  $\mathbf{y}_i - \mathbf{X}_i\hat{\beta} - \mathbf{Z}_i\hat{\mathbf{b}}_i$  to e.g. look at serial correlation. However,  $\hat{\mathbf{b}}_i$  very much depends on the normality assumption for  $\mathbf{b}$ , and is also influenced by the assumed structure for  $\mathbf{V}_i$ . Thus, these residuals are not well-suited to check assumptions on  $\mathbf{V}_i$ .

- **Influence diagnostics:** We can remove one subject from the analysis and recompute the parameter estimates. There are several distance measures to measure the influence of a single subject.

# Linear Mixed Models in R

See `examples_mixed_models.R` for example code for

- linear mixed models
- additive mixed models
- testing for variance components.

A reference for Linear Mixed Models using the `lme()` function in the `nlme` package in R: [Pinheiro & Bates, 2000](#).

A reference for Additive Mixed Models or Generalized Additive Mixed Models using the `mgcv` package in R: [Wood, 2006](#).

For many examples of Linear Mixed Models in SAS, see [Verbeke & Molenberghs, 2000](#).