

Stat/Biostat 571 Statistical Methodology: Regression Models for Dependent Data

Jon Wakefield

Departments of Statistics and Biostatistics, UW

Lectures: Monday/Wednesday/Friday 1.30–2.20, T473.

Coursework: (and approximate percentage contribution to final grade) weekly (30%). Examination at mid-term (30%) and final (40%). The mid-term and final will be takehome.

Office Hours:

Jon: Monday 2.30–3.20 and Wednesday 2.30–3.30 (Biostatistics, Health Sciences, 616-6292). Or by appointment (jonno@u.washington.edu, Padelford: 616-9388, HS: 616-6292).

TA: Youyi Fong (yfong@u); office hours to be arranged.

STAT/BIOSTAT 578 Data Analysis, strongly recommended for Applied Exam (Biostat students). 571 describes methods and not data analysis.

Computing will be carried out using R and WinBUGS.

Class website: <http://courses.washington.edu/b571/>

1

Textbooks:

Main Texts

Diggle, P.J., Heagerty, P., Liang, K.-Y. and Zeger, S.L. (2002). *Analysis of Longitudinal Data, Second Edition*. Oxford University Press.

Fitzmaurice, G.M., Laird, N.M. and Ware, J.H. (2004). *Applied Longitudinal Analysis*, Wiley.

Background Texts

Gelman, A., Carlin, J.B., Stern, H.S. and Rubin, D.B. (1995). *Bayesian Data Analysis*, CRC Press.

Hand, D. and Crowder, M.J. (1996). *Practical Longitudinal Data Analysis*, CRC Press.

Pinheiro, J. and Bates, D.G. (2000). *Mixed-Effects Models in S and S-PLUS*, Springer-Verlag,

Verbeke, G. and Molenberghs, G. (2000). *Linear Mixed Models for Longitudinal Data*. Springer-Verlag.

Davison, A.C. (2003). *Statistical Models*. Cambridge University Press.

Demidenko, E. (2004). *Mixed Models: Theory and Applications*, Wiley.

McCullagh, P. and Nelder, J.A. (1989) *Generalized Linear Models, Second Edition*, CRC Press.

2

COURSE OUTLINE

Revision

Motivating Datasets; Benefits and Challenges of Dependent Data; Marginal versus Conditional Modeling. Sandwich Estimation; Ordinary and Weighted Least Squares. Likelihood and Bayesian approaches.

Linear Models

Linear Mixed Effects Models; Frequentist and Bayesian Inference; Equivalence of Marginal and Conditional Modeling.

General Regression Models

Generalized Linear Mixed Models; Frequentist and Bayesian Inference; Non-equivalence of Marginal and Conditional Modeling.

Binary Data Models

Modeling the covariance structure. Mixed Effects approach.

Depending on Time:

Model Selection/Formulation

Spatial Models

3

Learning Objectives: 571, Winter 2008

By the end of the course the student should be able to:

- Explain why conventional regression models (as reviewed in 570) are inappropriate for dependent data, and explain the likely detrimental effects of their use in the dependent context.
- Describe the key differences between the generalized estimating equations, likelihood and Bayesian approaches to modeling dependent data. In particular the student should be able to describe the differences between the conditional and marginal approaches to modeling, and the advantages and drawbacks of each.
- Suggest appropriate approach(es) to modeling dependent outcomes, critically evaluate the fit of the fitted models, and interpret the estimated regression coefficients.
- Outline implementation strategies for each of the generalized estimating equations, likelihood and Bayesian approaches.

4

OVERVIEW

Recall: in a *regression analysis* we model a response, Y , as a function of covariates, \mathbf{x} .

In 570 we considered situations in which responses are *conditionally independent*, that is

$$\begin{aligned} p(Y_1, \dots, Y_n | \boldsymbol{\beta}, \mathbf{x}) &= p(Y_1 | \boldsymbol{\beta}, \mathbf{x}_1) \times p(Y_2 | Y_1, \boldsymbol{\beta}, \mathbf{x}_2) \times \dots \times p(Y_n | Y_1, \dots, Y_{n-1}, \boldsymbol{\beta}, \mathbf{x}_n) \\ &= p(Y_1 | \boldsymbol{\beta}, \mathbf{x}_1) \times p(Y_2 | \boldsymbol{\beta}, \mathbf{x}_2) \times \dots \times p(Y_n | \boldsymbol{\beta}, \mathbf{x}_n) \end{aligned}$$

so that observations are independent *given* parameters $\boldsymbol{\beta}$ and covariates $\mathbf{x}_1, \dots, \mathbf{x}_n$.

In general, Y_1, \dots, Y_n are *never* independent. For example, suppose

$$E[Y_i | \mu, \sigma^2] = \mu, \quad \text{var}(Y_i | \mu, \sigma^2) = \sigma^2,$$

$i = 1, 2$ and $\text{cov}(Y_1, Y_2 | \mu, \sigma^2) = 0$. Then if we are told y_1 , this will change the way we think about y_2 so that $p(Y_2 | Y_1) \neq p(Y_2)$, and the observations are not independent, however $p(Y_2 | Y_1, \mu, \sigma^2) = p(Y_2 | \mu, \sigma^2)$, so that we have conditional independence.

5

Motivating Examples

We distinguish between dependence induced by missing covariates, and that due to contagion (for example, in an infectious disease context) – we will not consider the latter.

One theme of the course will be modeling *residual* dependence, i.e. after we have controlled for covariates.

The obvious situations in which we would expect dependence is in data collected over time or space (but lots of others possible, e.g. families).

Example 1: Growth data

Table 1 records dental measurements of the distance in millimeters from the center of the pituitary gland to the pterygo-maxillary fissure in 11 girls and 16 boys at the ages of 8, 10, 12 and 14 years.

Here we have an example of *repeated measures* or *longitudinal* data.

Figure 1 plots these data and we see that dental growth for each child increases in an approximately linear fashion.

One common aim of such studies is to identify the *within-individual* and *between-individual* sources of variability.

6

Girls	8	10	12	14
1	21	20	21.5	23
2	21	21.5	24	25.5
3	20.5	24	24.5	26
4	23.5	24.5	25	26.5
5	21.5	23	22.5	23.5
6	20	21	21	22.5
7	21.5	22.5	23	25
8	23	23	23.5	24
9	20	21	22	21.5
10	16.5	19	19	19.5
11	24.5	25	28	28
Boys	8	10	12	14
1	26	25	29	31
2	21.5	22.5	23	26.5
3	23	22.5	24	27.5
4	25.5	27.5	26.5	27
5	20	23.5	22.5	26
6	24.5	25.5	27	28.5
7	22	22	24.5	26.5
8	24	21.5	24.5	25.5
9	23	20.5	31	26
10	27.5	28	31	31.5
11	23	23	23.5	25
12	21.5	23.5	24	28
13	17	24.5	26	29.5
14	22.5	25.5	25.5	26
15	23	24.5	26	30
16	22	21.5	23.5	25

7

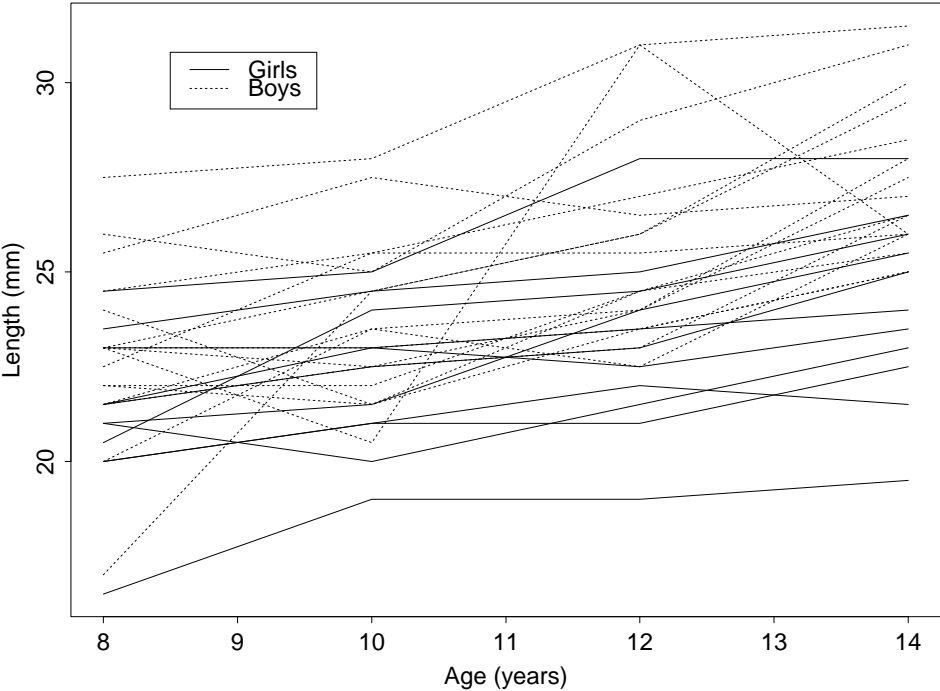


Figure 1: Dental growth data for girls and boys.

Inference

We may be interested in characterizing:

1. the *average* growth curve, or
2. the growth for a *particular* child.

Two types of analysis that will be distinguished are *marginal* and *conditional*. The former is designed for questions of type 1, and the latter may be used for both types, but requires more assumptions.

Even if the question of interest is of type 1, we still have to acknowledge the dependence of responses on the same individual – we do not have 11×4 independent observations on girls and 16×4 independent observations on boys but rather 11 and 16 *sets* of observations on girls and boys.

For either question of interest ignoring the dependence leads to incorrect standard errors and confidence/credible interval coverage.

A **marginal** approach (via Generalized Estimating Equations, GEE) to modeling specifies the moments of the data only, while in a **conditional** (via a Mixed Effects Model, MEM) approach the responses of specific individuals are modeled.

9

Models

First question is: why not just analyze the data from each child separately? Possible but we wouldn't be able to make formal statements about:

- The average growth rate of teeth for a girl in the age range 8–14 years.
- The between-girl variability in growth rates.

The totality of data on girls may also aid in the estimation of the growth rate for a particular girl – becomes more critical as the number of observations per child decreases. For example, in an extreme case, suppose a particular girl has only one measurement.

At the other extreme we could fit a single curve to the data from all of the girl's data together. The problem with this is that we do not have independent observations, and what if we are interested in inference for a particular child, or for a future child?

Example 2: Spatial Data

Dependent data may result from studies with a significant spatial component.

Split Plot Data

Example: Three varieties of oats, four nitrogen concentrations.

Agricultural land was grouped into six blocks, each with three plots, and with each plot further sub-divided into four sub-plots. Within each subplot a combination of oats and nitrogen was planted. Hence we have $6 \times 3 \times 4 = 72$ observations.

We would expect observations within the same block to be correlated.

11

Revision Material: Estimating Functions

Let $\mathbf{Y} = (Y_1, \dots, Y_n)$, represent n observations from a distribution indexed by a p -dimensional parameter $\boldsymbol{\theta}$, with $\text{cov}(Y_i, Y_j \mid \boldsymbol{\theta}) = 0$, $i \neq j$.

In the following, for ease of presentation, we assume that Y_i , $i = 1, \dots, n$ are independent and identically distributed (i.i.d.).

An *estimating function* is a function

$$\mathbf{G}_n(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n \mathbf{G}(\boldsymbol{\theta}, Y_i) \quad (1)$$

of the same dimension as $\boldsymbol{\theta}$ for which

$$\mathbf{E}[\mathbf{G}_n(\boldsymbol{\theta})] = \mathbf{0} \quad (2)$$

for all $\boldsymbol{\theta}$. The estimating function $\mathbf{G}_n(\boldsymbol{\theta})$ is a random variable because it is a function of \mathbf{Y} .

The corresponding *estimating equation* that defines the estimator $\hat{\boldsymbol{\theta}}_n$ has the form

$$\mathbf{G}_n(\hat{\boldsymbol{\theta}}_n) = \frac{1}{n} \sum_{i=1}^n \mathbf{G}(\hat{\boldsymbol{\theta}}_n, Y_i) = \mathbf{0}. \quad (3)$$

12

Result: Suppose that $\hat{\boldsymbol{\theta}}_n$ is a solution to the estimating equation

$$\mathbf{G}_n(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n \mathbf{G}(\boldsymbol{\theta}, Y_i) = \mathbf{0},$$

i.e. $\mathbf{G}_n(\hat{\boldsymbol{\theta}}_n) = \mathbf{0}$. Then $\hat{\boldsymbol{\theta}}_n \rightarrow_p \boldsymbol{\theta}$ (consistency) and

$$\sqrt{n} (\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}) \rightarrow_d N_p(\mathbf{0}, \mathbf{A}^{-1} \mathbf{B} \mathbf{A}^{\text{T}-1}) \quad (4)$$

(asymptotic normality) where

$$\mathbf{A} = \mathbf{A}(\boldsymbol{\theta}) = \text{E} \left[\frac{\partial}{\partial \boldsymbol{\theta}} \mathbf{G}(\boldsymbol{\theta}, Y) \right]$$

and

$$\mathbf{B} = \mathbf{B}(\boldsymbol{\theta}) = \text{E}[\mathbf{G}(\boldsymbol{\theta}, Y) \mathbf{G}(\boldsymbol{\theta}, Y)^{\text{T}}] = \text{cov}\{\mathbf{G}(\boldsymbol{\theta}, Y)\}.$$

The form of the variance in (4) has lead to it being named a **sandwich estimator**.

13

Example: Least Squares Estimation

For the ordinary least squares/maximum likelihood estimator

$$\hat{\boldsymbol{\beta}} = (\mathbf{x}^{\text{T}} \mathbf{x})^{-1} \mathbf{x}^{\text{T}} \mathbf{Y} \quad (5)$$

with

$$\text{var}(\hat{\boldsymbol{\beta}}) = (\mathbf{x}^{\text{T}} \mathbf{x})^{-1} \sigma^2$$

if $\text{var}(\mathbf{Y} \mid \mathbf{x}) = \sigma^2 \mathbf{I}$.

Suppose that $\text{var}(\mathbf{Y} \mid \mathbf{x}) = \sigma^2 \mathbf{V}$ so that the model from which the estimator (5) was derived was incorrect.

Then this estimator is still unbiased but the appropriate variance estimator is

$$\begin{aligned} \text{var}(\hat{\boldsymbol{\beta}}) &= (\mathbf{x}^{\text{T}} \mathbf{x})^{-1} \mathbf{x}^{\text{T}} \text{var}(\mathbf{Y} \mid \mathbf{x}) \mathbf{x} (\mathbf{x}^{\text{T}} \mathbf{x})^{-1} \\ &= (\mathbf{x}^{\text{T}} \mathbf{x})^{-1} \mathbf{x}^{\text{T}} \mathbf{V} \mathbf{x} (\mathbf{x}^{\text{T}} \mathbf{x})^{-1} \sigma^2 \end{aligned} \quad (6)$$

14

Expression (6) can also be derived directly from the estimating function

$$\mathbf{G}_n(\boldsymbol{\beta}) = \mathbf{x}^T(\mathbf{Y} - \mathbf{x}\boldsymbol{\beta}),$$

from which we know that

$$(\mathbf{A}_n^{-1} \mathbf{B}_n \mathbf{A}_n^T)^{1/2}(\hat{\boldsymbol{\beta}}_n - \boldsymbol{\beta}) \rightarrow_d N_{k+1}(\mathbf{0}, \mathbf{I}),$$

(note not iid observations here) where

$$\mathbf{B}_n = \text{var}(\mathbf{G}) = \mathbf{x}^T \mathbf{V} \mathbf{x} \sigma^2$$

and

$$\mathbf{A}_n = \text{E} \left[\frac{\partial \mathbf{G}}{\partial \boldsymbol{\beta}} \right] = -\mathbf{x}^T \mathbf{x},$$

to give

$$\text{var}(\hat{\boldsymbol{\beta}}) = (\mathbf{x}^T \mathbf{x})^{-1} \mathbf{x}^T \mathbf{V} \mathbf{x} (\mathbf{x}^T \mathbf{x})^{-1} \sigma^2.$$

We still need to know \mathbf{V} though.

15

Sandwich estimator with uncorrelated errors

We relax the constant variance assumptions. Consider the estimating function

$$\mathbf{G}(\boldsymbol{\beta}) = \mathbf{x}^T(\mathbf{Y} - \mathbf{x}\boldsymbol{\beta}).$$

The “bread” of the sandwich, \mathbf{A}^{-1} , remains unchanged since \mathbf{A} does not depend on \mathbf{Y} .

The “filling” becomes

$$\mathbf{B} = \text{var}\{\mathbf{G}\} = \mathbf{x}^T \text{var}(\mathbf{Y}) \mathbf{x} = \sum_{i=1}^n \sigma_i^2 \mathbf{x}_i^T \mathbf{x}_i, \quad (7)$$

where $\sigma_i^2 = \text{var}(Y_i)$ and we have assumed that the data are uncorrelated.

Unfortunately σ_i^2 is unknown – we now discuss various estimation methods.

An obvious estimator is given by

$$\widehat{\mathbf{B}}_n = \sum_{i=1}^n \mathbf{x}_i^T \mathbf{x}_i (Y_i - \mathbf{x}_i \widehat{\boldsymbol{\beta}})^2, \quad (8)$$

and its use provides a consistent estimator of (7), if the data are uncorrelated.

For linear regression the estimator

$$\widehat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (Y_i - \mathbf{x}_i \widehat{\boldsymbol{\beta}})^2 = \frac{1}{n} \sum_{i=1}^n \widehat{\sigma}_i^2,$$

is downwardly biased, with bias $-p\sigma^2/n$.

The sandwich estimator is therefore also downwardly biased.

Using

$$\widetilde{\sigma}_i^2 = \frac{n}{n-p} (Y_i - \mathbf{x}_i \widehat{\boldsymbol{\beta}})^2 \quad (9)$$

provides a simple correction, but in general the estimator of the variance has finite bias since the bias in $\widehat{\sigma}^2$ changes as a function of the design points \mathbf{x}_i – various corrections have been suggestions (see Kauermann and Carroll, 2001, *JASA*).

17

Likelihood Methods

A special case of the estimating function methodology occurs when the estimating equation

$$\mathbf{G} = \frac{\partial l}{\partial \boldsymbol{\theta}}$$

is a score equation (derivative of the log-likelihood). Then $\widehat{\boldsymbol{\theta}}$ is the MLE and

$$\sqrt{n} (\widehat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}) \rightarrow_d N_p(\mathbf{0}, \mathbf{I}^{-1}) \quad (10)$$

(asymptotic normality) where \mathbf{I} is the expected information matrix:

$$\mathbf{I} = \mathbf{A}(\boldsymbol{\theta}) = \mathbf{E} \left[\frac{\partial}{\partial \boldsymbol{\theta}} \mathbf{G}(\boldsymbol{\theta}, Y) \right] = \mathbf{B}(\boldsymbol{\theta}) = \mathbf{E}[\mathbf{G}(\boldsymbol{\theta}, Y) \mathbf{G}(\boldsymbol{\theta}, Y)^T] = \text{cov}\{\mathbf{G}(\boldsymbol{\theta}, Y)\}.$$

Bayesian Inference

In the Bayesian approach to inference all *unknown* quantities contained in a probability model for the observed data are treated as random variables.

These unknowns may include, for example, missing data, the true covariate value in an errors-in-variables setting, or the failure time of a censored survival observation.

Inference is made through the *posterior* probability distribution of $\boldsymbol{\theta}$ after observing \mathbf{y} , and is determined from Bayes theorem:

$$p(\boldsymbol{\theta} \mid \mathbf{y}) = \frac{p(\mathbf{y} \mid \boldsymbol{\theta}) \times \pi(\boldsymbol{\theta})}{p(\mathbf{y})},$$

where, for continuous $\boldsymbol{\theta}$, the normalizing constant is given by

$$p(\mathbf{y}) = \int_{\boldsymbol{\theta}} p(\mathbf{y} \mid \boldsymbol{\theta}) p(\boldsymbol{\theta}) \, d\boldsymbol{\theta},$$

and is the marginal probability of the observed data given the model (likelihood and prior). Ignoring this constant gives

$$\begin{aligned} p(\boldsymbol{\theta} \mid \mathbf{y}) &\propto p(\mathbf{y} \mid \boldsymbol{\theta}) \times \pi(\boldsymbol{\theta}) \\ \text{Posterior} &\propto \text{Likelihood} \times \text{Prior} \end{aligned}$$

19

The use of the posterior distribution for inference is very intuitively appealing since it probabilistically combines information on the parameters arising from the data and from prior beliefs.

An important observation is that for all $\boldsymbol{\theta}$ for which $\pi(\boldsymbol{\theta}) = 0$ we have $p(\boldsymbol{\theta} \mid \mathbf{y}) = 0$ also, regardless of any realization of the observed data. This has important consequences for prior specification and clearly shows that great care should be taken in excluding parts of the parameter space *a priori*.

Inference

To summarize the typically multivariate posterior distribution, $p(\boldsymbol{\theta} \mid \mathbf{y})$, marginal distributions for parameters of interest may be considered.

For example the univariate marginal distribution for a component θ_i is given by

$$p(\theta_i \mid \mathbf{y}) = \int_{\boldsymbol{\theta}_{-i}} p(\boldsymbol{\theta} \mid \mathbf{y}) \, d\boldsymbol{\theta}_{-i}, \quad (11)$$

where $\boldsymbol{\theta}_{-i}$ is the vector $\boldsymbol{\theta}$ excluding θ_i .

Posterior moments may be evaluated from the marginal distributions; for example the posterior mean is given by

$$E[\theta_i \mid \mathbf{y}] = \int_{\theta_i} \theta_i p(\theta_i \mid \mathbf{y}) \, d\theta_i. \quad (12)$$

Further summarization may be carried out to yield the $100 \times q\%$ quantile, $\theta_i(q)$ ($0 < q < 1$) by solving

$$\int_{-\infty}^{\theta_i(q)} p(\theta_i \mid \mathbf{y}) \, d\theta_i. \quad (13)$$

In particular, the posterior median, $\theta_i(0.5)$, will often provide an adequate summary of the location of the posterior marginal distribution.

21

A $100 \times p\%$ equi-tailed *credible interval* ($0 < p < 1$) is provided by $[\theta_i\{(1-p)/2\}, \theta_i\{(1+p)/2\}]$.

Such intervals are usually reported though in some cases it which the posterior is skewed one may wish to instead calculate a *highest posterior density* (HPD) interval in which points inside the interval have higher posterior density than those outside the interval (such an interval is also the shortest credible interval).

Another useful inferential quantity is the *predictive* distributions for future observations \mathbf{z} which is given, under conditional independence, by

$$p(\mathbf{z} \mid \mathbf{y}) = \int_{\boldsymbol{\theta}} p(\mathbf{z} \mid \boldsymbol{\theta}) p(\boldsymbol{\theta} \mid \mathbf{y}) \, d\boldsymbol{\theta}. \quad (14)$$

This clearly assumes that the system under study is stable so that the likelihood for future observations is still the relevant data generation mechanism.

Bayesian inference is deceptively simple to describe probabilistically, but there have been two major obstacles to its routine use. The first is how to specify prior distributions and the second is how to evaluate the integrals required for inference, for example, (11)–(14), given that for most models, these are analytically intractable

Example: Normal linear regression, variance unknown

Suppose we have $Y_i \mid \boldsymbol{\beta}, \sigma^2 \sim_{ind} N(\mathbf{x}_i \boldsymbol{\beta}, \sigma^2)$, $i = 1, \dots, n$. $\dim(\boldsymbol{\beta}) = p$.

MLE: $\hat{\boldsymbol{\beta}} \sim t_p(\boldsymbol{\beta}, (\mathbf{x}^T \mathbf{x})^{-1} s^2, n - p)$, a Student t distribution with $n - p$ degrees of freedom.

Improper prior: $\pi(\boldsymbol{\beta}, \sigma^2) \propto \sigma^{-2}$.

Marginal posterior:

$$p(\boldsymbol{\beta} \mid \mathbf{y}) = \int p(\boldsymbol{\beta}, \sigma^2 \mid \mathbf{y}) d\sigma^2,$$

where

$$p(\boldsymbol{\beta}, \sigma^2 \mid \mathbf{y}) \propto l(\boldsymbol{\beta}, \sigma^2) \times \pi(\boldsymbol{\beta}, \sigma^2).$$

Hence

$$\begin{aligned} p(\boldsymbol{\beta} \mid \mathbf{y}) &= \int \frac{(2\pi\sigma^2)^{-n/2}}{\sigma^2} \exp \left\{ -\frac{[(n-p)s^2 + (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})^T \mathbf{x}^T \mathbf{x} (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})]}{2\sigma^2} \right\} d\sigma^2 \\ &\propto \int (\sigma^2)^{-(n/2+1)} \exp \left\{ -\frac{c}{2\sigma^2} \right\} d\sigma^2 \end{aligned}$$

where

$$c = (n-p)s^2 + (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})^T \mathbf{x}^T \mathbf{x} (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}).$$

23

We have the kernel of an inverse Gamma distribution $\text{IGa}(n/2, c)$.

An inverse gamma r.v. X has density

$$p(x) = \frac{\beta^\alpha}{\Gamma(\alpha)} x^{-(\alpha+1)} \exp(-\beta/x), \quad x > 0.$$

Hence

$$\begin{aligned} p(\boldsymbol{\beta} \mid \mathbf{y}) &\propto \left(\frac{c}{2}\right)^{-n/2} \\ &\propto \{(n-p)s^2 + (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})^T \mathbf{x}^T \mathbf{x} (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})\}^{-n/2} \\ &\propto \left\{ 1 + \frac{(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})^T \mathbf{x}^T \mathbf{x} (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})}{(n-p)s^2} \right\}^{[-(n-p)+p]/2} \\ &= \left\{ 1 + \frac{(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})^T \Sigma^{-1} (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})}{n-p} \right\}^{[-(n-p)+p]/2} \end{aligned}$$

where $\Sigma = (\mathbf{x}^T \mathbf{x})^{-1} s^2$.

24

Hence the posterior

$$\boldsymbol{\beta} \mid \mathbf{y} \sim t_p(\hat{\boldsymbol{\beta}}, (\mathbf{x}^T \mathbf{x})^{-1} s^2, n - p).$$

A p dimensional multivariate Student's t r.v. \mathbf{X} with degrees of freedom d has density

$$p(\mathbf{x}) = \frac{\Gamma\{(d+p)/2\}}{\Gamma(d/2)(d\pi)^{p/2}} |\boldsymbol{\Sigma}|^{-1/2} \times [1 + (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})/d]^{-(d+p)/2}.$$

Note the similarity with frequentist inference.

25

LINEAR MODELS

We now begin thinking about specific situations, starting with linear models. Clearly, in general, ignoring dependence will give inappropriate standard errors.

While making inference for dependent data is more difficult than for independent data, designs that collect dependent data can be very efficient. For example (as we see shortly), in a longitudinal data setting applying different treatments to the same patient over time can be very beneficial since each patient acts as their own control.

While in the Bayesian approach to inference all parameters are viewed as random variables, in the frequentist approach there is a distinction between *fixed effects* (unknown constants) and *random effects* (random variables from a distribution).

26

Design Implications of a Longitudinal Study

To examine the implications of carrying out a longitudinal study, as compared to a cross-sectional study, we consider a very simple situation in which we wish to compare two treatments, coded as -1 and +1, and we have a linear model.

Cross-Sectional Study:

A single measurement is taken on each of $m = 4$ individuals where

$$Y_{i1} = \beta_0 + \beta_1 x_{i1} + \epsilon_{i1},$$

$i = 1, \dots, m = 4$, ϵ_{i1} iid with $\text{var}(\epsilon_{i1}) = \sigma^2$ and

$$x_{11} = -1, x_{21} = -1, x_{31} = 1, x_{41} = 1$$

Note: $E[Y_1|x = 1] - E[Y_1|x = -1] = 2\beta_1$.

Using least squares:

$$\hat{\beta}_0^c = \frac{\sum_{i=1}^4 Y_{i1}}{4}, \quad \hat{\beta}_1^c = \frac{Y_{31} + Y_{41} - (Y_{11} + Y_{21})}{4},$$

and

$$\text{var}(\hat{\beta}_0^c) = \text{var}(\hat{\beta}_1^c) = \frac{\sigma^2}{4}.$$

27

Longitudinal Study:

We assume the model

$$Y_{ij} = \beta_0 + \beta_1 x_{ij} + b_i + \delta_{ij},$$

with b_i and δ_{ij} independent and with $\text{var}(b_i) = \sigma_0^2$, $\text{var}(\delta_{ij}) = \sigma_\delta^2$. We therefore have marginally:

$$\text{var}(Y_{ij}|\beta_0, \beta_1) = \sigma_0^2 + \sigma_\delta^2 = \sigma^2,$$

and

$$\text{cov}(Y_{i1}, Y_{i2}) = \sigma_0^2.$$

We let $\rho = \sigma_0^2/\sigma^2$, represent the correlation on observations on the same individual.

We consider two situations, both with two observations on two individuals:

Constant treatment for each individual:

$$x_{11} = x_{12} = -1, \quad x_{21} = x_{22} = 1.$$

Changing treatment for each individual:

$$x_{11} = x_{22} = 1, \quad x_{12} = x_{21} = -1.$$

Using Generalized Least Squares we have

$$\hat{\boldsymbol{\beta}}^l = (\mathbf{x}^T \mathbf{R}^{-1} \mathbf{x})^{-1} \mathbf{x}^T \mathbf{R}^{-1} \mathbf{Y},$$

and

$$\text{var}(\hat{\boldsymbol{\beta}}^l) = (\mathbf{x}^T \mathbf{R}^{-1} \mathbf{x})^{-1} \sigma^2,$$

where

$$\mathbf{R} = \begin{bmatrix} 1 & \rho & 0 & 0 \\ \rho & 1 & 0 & 0 \\ 0 & 0 & 1 & \rho \\ 0 & 0 & \rho & 1 \end{bmatrix}.$$

In lectures we will show that

$$\text{var}(\hat{\beta}_1^l) = \frac{\sigma^2(1 - \rho^2)}{4 - 2\rho(x_{11}x_{12} + x_{21}x_{22})}.$$

29

The *efficiency* e is given by

$$e = \frac{\text{var}(\hat{\beta}_1^l)}{\text{var}(\hat{\beta}_1^c)} = \frac{(1 - \rho^2)}{1 - \rho(x_{11}x_{12} + x_{21}x_{22})/2}.$$

Usually we have $\rho > 0$.

For the constant treatment longitudinal study

$$e = 1 + \rho,$$

so that the cross-sectional study is preferable since we have lost information due to the correlation.

For the changing treatment longitudinal study

$$e = 1 - \rho,$$

so that the longitudinal study is more efficient, because each individual is acting as their own control, that is, we are making within-individual comparisons.

If $\rho = 0$ the designs have the same efficiency.

Example: Dental Growth Data

Suppose $\hat{\beta}_0^m$ and $\hat{\beta}_1^m$ are the marginal intercept and slope estimates, and let

$$e_{ij}^m = Y_{ij} - \hat{\beta}_0^m - \hat{\beta}_1^m t_j,$$

$i = 1, \dots, 11; j = 1, \dots, 4$, denote marginal residuals, and

$$\begin{bmatrix} \sigma_1 & & & \\ \rho_{12} & \sigma_2 & & \\ \rho_{13} & \rho_{23} & \sigma_3 & \\ \rho_{14} & \rho_{24} & \rho_{34} & \sigma_4 \end{bmatrix} \quad (15)$$

represent the standard deviation/correlation matrix of the residuals, where

$$\sigma_j = \sqrt{\text{var}(e_{ij}^m)},$$

is the variance of the length at time t_j , $j = 1, \dots, 4$, and

$$\rho_{jk} = \frac{\text{cov}(e_{ij}^m, e_{ik}^m)}{\sqrt{\text{var}(e_{ij}^m)\text{var}(e_{ik}^m)}},$$

is the correlation between residual measurements at times t_j and t_k taken on the same girl, $j \neq k, j, k = 1, \dots, 4$.

31

Across girls we may empirically estimate the entries of (15) by

$$\begin{bmatrix} 2.12 & & & \\ 0.83 & 1.90 & & \\ 0.86 & 0.90 & 2.36 & \\ 0.84 & 0.88 & 0.95 & 2.44 \end{bmatrix} \quad (16)$$

illustrating that there is a suggestion that the variance is increasing with the mean, and clear correlation between residuals at different times on the same girl.

The fitting of a single curve, and using methods for independent data, ignores the correlations within each child's data and so standard errors will clearly be inappropriate.

Fitting a marginal model such as this is appealing in one sense, however, since it allows the direct comparison of the average responses in different (in this example the populations of girls at different ages) and forms the basis of the generalized estimating equations (GEE) approach

An alternative fixed effects approach is to assume a fixed curve for each child and analyze each set of data separately.

We will also often be interested in making formal inference for the population of girls from which the eleven in the data are viewed as a random sample. This forms the basis of the mixed effects model approach.

Figure 2(b) displays the lines corresponding to each of these fixed effects approaches.

33

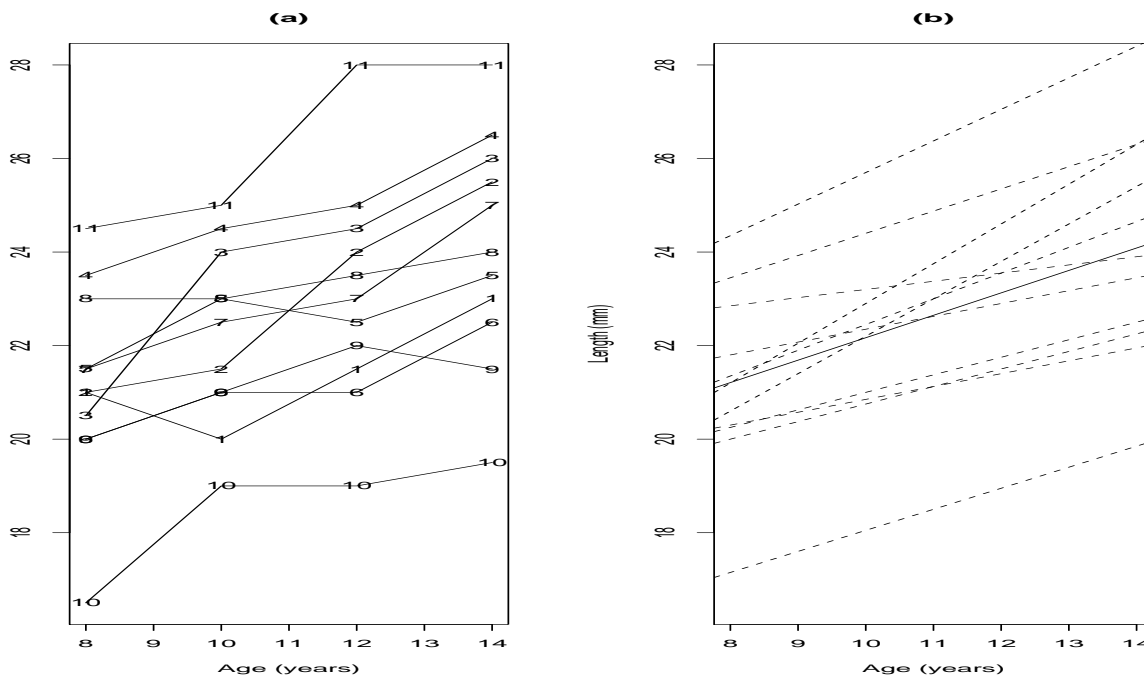


Figure 2: Dental plots for girls only: (a) Individual observed data (with plotting symbol girl index), (b) Individual fitted curves (dashed) and overall fitted curve (solid).

34

Linear Mixed Effects Models

The basic idea behind mixed effects models is to assume that each unit has a regression model characterized by unit-specific parameters, with these parameters being a combination of fixed effects that are common to all units in the population, and then unit-specific perturbations, or random effects (hence “mixed” effects refers to the combination of fixed and random effects).

Given data $\mathbf{y}_i = (y_{i1}, \dots, y_{in_i})^T$ on unit i a mixed effects model is characterized by a combination of

- a $(k + 1) \times 1$ vector of fixed effects, $\boldsymbol{\beta}$,
- a $(q + 1) \times 1$ vector of random effects, \mathbf{b}_i , with $q \leq k$.
- $\mathbf{x}_i = (\mathbf{x}_{i1}, \dots, \mathbf{x}_{in_i})^T$, the design matrix for the fixed effect with $\mathbf{x}_{ij} = (1, x_{ij1}, \dots, x_{ijk})^T$, and
- $\mathbf{z}_i = (\mathbf{z}_{i1}, \dots, \mathbf{z}_{in_i})^T$, and design matrix for the random effects with $\mathbf{z}_{ij} = (1, z_{ij1}, \dots, z_{ijq})^T$.

35

We then have the following (two stage) Linear Mixed Effects Model (LMEM):

Stage 1: Response model, *conditional* on random effects:

$$\mathbf{y}_i = \mathbf{x}_i \boldsymbol{\beta} + \mathbf{z}_i \mathbf{b}_i + \boldsymbol{\epsilon}_i, \quad (17)$$

where $\boldsymbol{\epsilon}_i$ is an $n_i \times 1$ zero mean vector of error terms.

Stage 2: Model for random terms:

$$\begin{aligned} E[\boldsymbol{\epsilon}_i] &= \mathbf{0}, \quad \text{var}(\boldsymbol{\epsilon}_i) = \mathbf{E}_i(\boldsymbol{\alpha}), \\ E[\mathbf{b}_i] &= \mathbf{0}, \quad \text{var}(\mathbf{b}_i) = \mathbf{D}(\boldsymbol{\alpha}), \\ \text{cov}(\mathbf{b}_i, \boldsymbol{\epsilon}_i) &= \mathbf{0} \end{aligned}$$

where $\boldsymbol{\alpha}$ is the vector of variance-covariance parameters.

The two stages define the marginal model:

$$\begin{aligned} E[\mathbf{y}_i] &= \boldsymbol{\mu}_i(\boldsymbol{\beta}) = \mathbf{x}_i \boldsymbol{\beta}, \\ \text{var}(\mathbf{y}_i) &= \mathbf{V}_i(\boldsymbol{\alpha}) = \mathbf{z}_i \mathbf{D} \mathbf{z}_i^T + \mathbf{E}_i, \\ \text{cov}(\mathbf{y}_i, \mathbf{y}_{i'}) &= \mathbf{0}, \quad i \neq i'. \end{aligned}$$

We describe likelihood and Bayesian approaches to inference.

36

Likelihood Inference

We need to specify a complete probability distribution for the data, and this follows by specifying distributions for $\boldsymbol{\epsilon}_i$ and \mathbf{b}_i , $i = 1, \dots, m$. A common model is

$$\boldsymbol{\epsilon}_i \sim_{ind} N(\mathbf{0}, \sigma_\epsilon^2 \mathbf{I}_{n_i}), \quad \mathbf{b}_i \sim_{iid} N(\mathbf{0}, \mathbf{D}),$$

where

$$\mathbf{D} = \begin{bmatrix} \sigma_{00}^2 & \sigma_{01}^2 & \cdots & \sigma_{0q}^2 \\ \sigma_{10}^2 & \sigma_{11}^2 & \cdots & \sigma_{1q}^2 \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{q0}^2 & \sigma_{q1}^2 & \cdots & \sigma_{qq}^2 \end{bmatrix}.$$

Here $\boldsymbol{\alpha} = (\sigma_\epsilon^2, \mathbf{D})$ denote the variance-covariance parameters. Here $\mathbf{V} = \mathbf{zDz}^T + \sigma_\epsilon^2 \mathbf{I}_N$, where $N = \sum_{i=1}^m n_i$.

Likelihood methods are designed for fixed effects, and so we integrate the random effects from the two-stage model:

$$p(\mathbf{y}|\boldsymbol{\beta}, \boldsymbol{\alpha}) = \int_{\mathbf{b}} p(\mathbf{y}|\mathbf{b}, \boldsymbol{\beta}, \boldsymbol{\alpha}) \times p(\mathbf{b}|\boldsymbol{\beta}, \boldsymbol{\alpha}) d\mathbf{b}.$$

37

Exploiting conditional independencies we have:

$$p(\mathbf{y}|\boldsymbol{\beta}, \boldsymbol{\alpha}) = \prod_{i=1}^m \int_{\mathbf{b}_i} p(\mathbf{y}_i|\mathbf{b}_i, \boldsymbol{\beta}, \sigma_\epsilon^2) \times p(\mathbf{b}_i|\mathbf{D}) d\mathbf{b}_i.$$

Since a convolution of normals is normal we obtain

$$\mathbf{y}|\boldsymbol{\beta}, \boldsymbol{\alpha} \sim \prod_{i=1}^m N\{\boldsymbol{\mu}_i(\boldsymbol{\beta}), \mathbf{V}_i(\boldsymbol{\alpha})\}.$$

The log-likelihood is

$$\begin{aligned} l(\boldsymbol{\beta}, \boldsymbol{\alpha}) = & - \frac{N}{2} \log 2\pi - \frac{1}{2} \sum_{i=1}^m \log |\mathbf{V}_i(\boldsymbol{\alpha})| \\ & - \frac{1}{2} \sum_{i=1}^m (\mathbf{Y}_i - \mathbf{x}_i \boldsymbol{\beta})^T \mathbf{V}(\boldsymbol{\alpha})_i^{-1} (\mathbf{Y}_i - \mathbf{x}_i \boldsymbol{\beta}). \end{aligned} \quad (18)$$

Example: One-way ANOVA

Consider the simple ANOVA model

$$Y_{ij} = \beta_0 + b_i + \epsilon_{ij},$$

with b_i and ϵ_{ij} independent and distributed as

- $b_i \sim_{ind} N(0, \sigma_0^2)$,
- $\epsilon_{ij} \sim_{ind} N(0, \sigma_\epsilon^2)$

for $i = 1, \dots, m$, $j = 1, \dots, n_i$, with $\sum_{i=1}^m n_i = N$. This model can also be written as

$$\mathbf{Y}_i = \mathbf{1}_n \beta_0 + \mathbf{1}_n b_i + \boldsymbol{\epsilon}_i,$$

with $E[\mathbf{Y}] = \mathbf{1}_N \beta_0$, $\text{var}(\mathbf{Y}) = \mathbf{V} = \mathbf{1}_N \mathbf{1}_N^T \sigma_0^2 + \mathbf{I}_N \sigma_\epsilon^2 = \mathbf{J}_N \sigma_0^2 + \mathbf{I}_N \sigma_\epsilon^2$, where \mathbf{J}_N is the $N \times N$ matrix of 1's.

39

The marginal variance \mathbf{V} is the $N \times N$ matrix

$$\sigma^2 \begin{bmatrix} 1 & \rho & \rho & \rho & \dots & 0 & 0 & 0 & 0 \\ \rho & 1 & \rho & \rho & \dots & 0 & 0 & 0 & 0 \\ \rho & \rho & 1 & \rho & \dots & 0 & 0 & 0 & 0 \\ \rho & \rho & \rho & 1 & \dots & 0 & 0 & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots & \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & 0 & 0 & \dots & 1 & \rho & \rho & \rho \\ 0 & 0 & 0 & 0 & \dots & \rho & 1 & \rho & \rho \\ 0 & 0 & 0 & 0 & \dots & \rho & \rho & 1 & \rho \\ 0 & 0 & 0 & 0 & \dots & \rho & \rho & \rho & 1 \end{bmatrix}$$

with $\sigma^2 = \sigma_\epsilon^2 + \sigma_0^2$ and

$$\rho = \frac{\sigma_0^2}{\sigma^2} = \frac{\sigma_0^2}{\sigma_\epsilon^2 + \sigma_0^2}.$$

40

Here we have a total of 3 regression parameters and variance components $(\beta_0, \sigma_0^2, \sigma_\epsilon^2)$, but $m + 3$ if we count the random effects.

A fixed effects model with a separate parameter for each group would have $m + 1$ parameters (and corresponds to the above model with $\sigma_0^2 = \infty$).

In some situations we may have more fixed and random effects than data points, but the random effects have a special status, since they are tied together through a common distribution.

Random effects may be viewed as a means by which dependencies are induced in marginal models.

41

Inference for Regression Parameters

The score equation for β is

$$\frac{\partial l}{\partial \beta} = \sum_{i=1}^m \mathbf{x}_i^T \mathbf{V}_i^{-1} \mathbf{Y}_i - \sum_{i=1}^m \mathbf{x}_i^T \mathbf{V}_i^{-1} \mathbf{x}_i \beta,$$

and yields the MLE for β as

$$\hat{\beta} = \left(\sum_{i=1}^m \mathbf{x}_i^T \mathbf{V}_i^{-1} \mathbf{x}_i \right)^{-1} \left(\sum_{i=1}^m \mathbf{x}_i^T \mathbf{V}_i^{-1} \mathbf{y}_i \right), \quad (19)$$

which is a weighted least squares estimator. If $\mathbf{D} = \mathbf{0}$ then $\mathbf{V} = \sigma_\epsilon^2 \mathbf{I}_N$ and $\hat{\beta}$ corresponds to the ordinary least squares estimator.

The variance of $\hat{\beta}$ may be obtained either directly from (19), or from the second derivative of the log-likelihood. Since

$$\frac{\partial^2 l}{\partial \beta \partial \beta^T} = - \sum_{i=1}^m \mathbf{x}_i^T \mathbf{V}_i^{-1} \mathbf{x}_i,$$

the observed and expected information matrices coincide with

$$\mathbf{I}_{\beta\beta} = -\mathbf{E} \left[\frac{\partial^2 l}{\partial \beta \partial \beta^T} \right] = \sum_{i=1}^m \mathbf{x}_i^T \mathbf{V}_i^{-1} \mathbf{x}_i.$$

42

The estimator, $\hat{\beta}$ is a linear combination of \mathbf{Y}_i and so, under correct specification of the model $\hat{\beta}$ is linear also and

$$\hat{\beta} \sim N_{k+1} \left\{ \beta, \left(\sum_{i=1}^m \mathbf{x}_i \mathbf{V}_i^{-1} \mathbf{x}_i \right)^{-1} \right\}.$$

In practice, α is never known, but asymptotically, as $m \rightarrow \infty$ (it is not sufficient to have m fixed and $n_i \rightarrow \infty$ for $i = 1, \dots, m$):

$$\left(\sum_{i=1}^m \mathbf{x}_i \mathbf{V}_i(\hat{\alpha})^{-1} \mathbf{x}_i \right)^{1/2} (\hat{\beta}_m - \beta) \rightarrow_d N_{k+1}(\mathbf{0}_{k+1}, \mathbf{I}_{k+1}),$$

where $\hat{\alpha}$ is a consistent estimator of α . This result is also relevant if the data and random effects are not normal, so long as the second moment assumptions are correct.

Various t and F -like approaches have been suggested for correcting for the estimation of α , see Verbeke and Molenberghs (2000, Chapter 6), but if the sampling size is not sufficiently large for reliable estimation of α , we recommend following a Bayesian approach to inference.

43

So far as the MLE is concerned, the expected information matrix is partitioned as

$$\mathbf{I}(\beta, \alpha) = \begin{bmatrix} \mathbf{I}_{\beta\beta} & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_{\alpha\alpha} \end{bmatrix}.$$

Standard ML theory gives the asymptotic distribution for the MLE $\hat{\beta}, \hat{\alpha}$, as

$$\begin{bmatrix} \hat{\beta} \\ \hat{\alpha} \end{bmatrix} \sim N_{k+1+r+1} \left(\begin{bmatrix} \beta \\ \alpha \end{bmatrix}, \begin{bmatrix} \mathbf{I}_{\beta\beta}^{-1} & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_{\alpha\alpha}^{-1} \end{bmatrix} \right),$$

where r is the number of distinct elements in \mathbf{D} .

We have already seen the form of $\mathbf{I}_{\beta\beta}$; the form of $\mathbf{I}_{\alpha\alpha}$ is not pleasant.

The diagonal form of the expected information has a number of implications. Firstly, we may carry out separate maximization of the log-likelihood with respect to β and α . Secondly, asymptotically, we have independence between $\hat{\beta}$ and $\hat{\alpha}$, so any consistent estimator of α will give an asymptotically efficient estimator for β .

Likelihood ratio tests are available for regression parameters.

Inference for Variance Components by MLE

The MLE of α follows from maximization of (18), and in general there is no closed-form solution.

The maximization may produce a negative variance estimate, in which case this variance is set equal to zero (MLEs must lie in the parameter space).

Maximum likelihood for variance components give estimators that do not acknowledge the estimation of β .

For the simple linear model, the MLE of σ^2 is RSS/n , and not the unbiased version $\text{RSS}/(n - k - 1)$.

An alternative and often preferable method is provided by restricted maximum likelihood (REML). In general REML is carried out to reduce bias in estimation of variance components by “accounting for the estimation of β ”.

45

Hypothesis tests for variance components

Testing whether random effect variances are zero requires care since the null hypothesis lies on the boundary, and so the usual regularity conditions are not satisfied.

As an example, in the model

$$Y_{ij} = \beta_0 + b_i + \mathbf{x}_{ij}\beta + \epsilon_{ij}$$

with $b_i \sim N(0, \sigma_0^2)$, consider the test of $H_0 : \sigma_0^2 = 0$ versus $H_A : \sigma_0^2 > 0$, where σ_0^2 is a non-negative scalar. In this case the asymptotic null distribution is a 50:50 mixture of χ_0^2 and χ_1^2 distributions, where the former is the distribution that gives probability mass 1 to the value 0.

Intuition: Estimating σ_0^2 is equivalent to estimating $\rho = \sigma_0^2/\sigma^2$, and setting equal to zero if the estimated correlation is negative, and under the null this will happen half the time. Setting $\hat{\rho} = 0$ gives the null, and so the likelihood ratio will be one.

If the usual χ_1^2 distribution is used then the null would be accepted too often, leading to a variance component structure that is too simple.

46

Inference for Variance Components by REML

Restricted (or residual) maximum likelihood (REML) is a method that has been proposed as an alternative to ML, there are a number of justifications; we later provide a Bayesian justification, and here provide another based on marginal likelihood.

Marginal Likelihood

Let $\mathbf{S}_1, \mathbf{S}_2, \mathbf{A}$ be a minimal sufficient statistic where \mathbf{A} is ancillary, and for which

$$\begin{aligned} p(\mathbf{y} \mid \boldsymbol{\lambda}, \boldsymbol{\phi}) &\propto p(\mathbf{s}_1, \mathbf{s}_2, \mathbf{a} \mid \boldsymbol{\lambda}, \boldsymbol{\phi}) \\ &= p(\mathbf{a})p(\mathbf{s}_1 \mid \mathbf{a}, \boldsymbol{\lambda})p(\mathbf{s}_2 \mid \mathbf{s}_1, \mathbf{a}, \boldsymbol{\lambda}, \boldsymbol{\phi}) \end{aligned}$$

where $\boldsymbol{\lambda}$ are parameters of interest and $\boldsymbol{\phi}$ are the remaining (nuisance) parameters.

47

Inference for $\boldsymbol{\lambda}$ may be based on the *marginal* likelihood

$$L_m(\boldsymbol{\lambda}) = p(\mathbf{s}_1 \mid \mathbf{a}, \boldsymbol{\lambda}).$$

This is desirable if inference is simplified or if it avoids problems encountered with standard likelihood methods. For example $\dim(\boldsymbol{\phi})$ may increase with n . The marginal likelihood has similar properties to a regular likelihood.

These advantages may outway the loss of efficiency in ignoring the $p(\mathbf{s}_2 \mid \mathbf{s}_1, \mathbf{a}, \boldsymbol{\lambda}, \boldsymbol{\phi})$ term. If there is no ancillary statistic then the marginal likelihood is

$$L_m(\boldsymbol{\lambda}) = p(\mathbf{s}_1 \mid \boldsymbol{\lambda}).$$

Example: Normal linear model

Assume $\mathbf{Y} \mid \boldsymbol{\beta}, \sigma^2 \sim_{ind} N_n(\mathbf{x}\boldsymbol{\beta}, \sigma^2 \mathbf{I})$ where $\dim(\boldsymbol{\beta}) = k + 1$. Suppose the parameter of interest is $\lambda = \sigma^2$, with remaining parameters $\boldsymbol{\phi} = \boldsymbol{\beta}$. Minimal sufficient statistics are: $s_1 = s^2 = \text{RSS}/(n - k - 1)$, and $\mathbf{s}_2 = \hat{\boldsymbol{\beta}}$. We have

$$p(\mathbf{y} \mid \sigma^2, \boldsymbol{\beta}) = p(s_1, \mathbf{s}_2 \mid \sigma^2, \boldsymbol{\beta}) = p(s_1 \mid \sigma^2)p(\mathbf{s}_2 \mid \boldsymbol{\beta}, \sigma^2).$$

Hence the marginal likelihood is

$$L_m(\sigma^2) = p(s^2 \mid \sigma^2).$$

We know

$$\frac{(n - k - 1)s^2}{\sigma^2} \sim \chi_{n-k-1}^2 = \text{Ga}\left(\frac{n - k - 1}{2}, \frac{1}{2}\right),$$

and so

$$p(s^2 \mid \sigma^2) = \left(\frac{n - k - 1}{2\sigma^2}\right)^{(n-k-1)/2} \frac{(s^2)^{(n-k-1)/2-1}}{\Gamma\left(\frac{n-k-1}{2}\right)} \times \exp\left[-\frac{(n - k - 1)s^2}{2\sigma^2}\right],$$

to give

$$l_m = \log L_m = -(n - k - 1) \log \sigma - \frac{(n - k - 1)s^2}{2\sigma^2},$$

and

$$\hat{\sigma}^2 = s^2.$$

49

REML for LMEM

To use marginal likelihood we need to find a function of the data, $\mathbf{U} = f(\mathbf{Y})$, whose distribution does not depend upon $\boldsymbol{\beta}$, and then base inference for $\boldsymbol{\alpha}$ on this distribution.

A natural function to choose is the vector of residuals following an ordinary least squares fit:

$$\begin{aligned} \mathbf{R} &= \mathbf{Y} - \mathbf{x}\hat{\boldsymbol{\beta}}_o = \mathbf{Y} - \mathbf{x}(\mathbf{x}^\top \mathbf{x})^{-1} \mathbf{x}^\top \mathbf{Y} \\ &= (\mathbf{I} - \mathbf{x}(\mathbf{x}^\top \mathbf{x})^{-1} \mathbf{x}^\top) \mathbf{Y} = (\mathbf{I} - \mathbf{H}) \mathbf{Y}, \end{aligned}$$

where $\hat{\boldsymbol{\beta}}_o = (\mathbf{x}^\top \mathbf{x})^{-1} \mathbf{x}^\top \mathbf{Y}$ is the OLS estimator.

We have

$$\mathbf{R} = (\mathbf{I} - \mathbf{H}) \mathbf{Y} = (\mathbf{I} - \mathbf{H})(\mathbf{x}\boldsymbol{\beta} + \mathbf{z}\mathbf{b} + \boldsymbol{\epsilon}) = (\mathbf{I} - \mathbf{H})(\mathbf{z}\mathbf{b} + \boldsymbol{\epsilon}),$$

and so the distribution of \mathbf{R} does not depend on $\boldsymbol{\beta}$.

Unfortunately the distribution of \mathbf{R} is degenerate as it has rank $N - k - 1$.

Consider the $(N - k - 1) \times 1$ random variables

$$\mathbf{U} = \mathbf{B}^T \mathbf{Y}$$

where \mathbf{B} is an $N \times (N - k - 1)$ matrix with $\mathbf{B}\mathbf{B}^T = \mathbf{I} - \mathbf{H}$ and $\mathbf{B}^T \mathbf{B} = \mathbf{I}$ (such a matrix always exists).

Then

$$\mathbf{U} = \mathbf{B}^T \mathbf{Y} = \mathbf{B}^T \mathbf{B} \mathbf{B}^T \mathbf{Y} = \mathbf{B}^T (\mathbf{I} - \mathbf{H}) \mathbf{Y} = \mathbf{B}^T \mathbf{R},$$

and $\mathbf{B}^T \mathbf{Y}$ is a linear combination of residuals.

Further $\mathbf{B}^T \mathbf{x} = \mathbf{0}$, so that

$$\mathbf{U} = \mathbf{B}^T \mathbf{Y} = \mathbf{B}^T \mathbf{z} \mathbf{b} + \mathbf{B}^T \boldsymbol{\epsilon},$$

and the distribution of \mathbf{U} does not depend upon $\boldsymbol{\beta}$, and $E[\mathbf{U}] = \mathbf{0}$.

We now derive the distribution of \mathbf{U} . To do this we consider the transformation from $\mathbf{Y} \rightarrow (\mathbf{U}, \hat{\boldsymbol{\beta}}_G) = (\mathbf{B}^T \mathbf{Y}, \mathbf{G}^T \mathbf{Y})$, where

$$\hat{\boldsymbol{\beta}}_G = \mathbf{G}^T \mathbf{Y} = (\mathbf{x}^T \mathbf{V}^{-1} \mathbf{x})^{-1} \mathbf{x}^T \mathbf{V}^{-1} \mathbf{Y},$$

is the generalized least squares estimator.

51

We derive the Jacobian of the transformation. To do this we need the following two facts:

1. $\det(\mathbf{A}^T \mathbf{A}) = \det(\mathbf{A}^T) \det(\mathbf{A}) = \det(\mathbf{A})^2$.
2. $\left| \begin{array}{cc} \mathbf{T} & \mathbf{U} \\ \mathbf{V} & \mathbf{W} \end{array} \right| = | \mathbf{T} || \mathbf{W} - \mathbf{V} \mathbf{T}^{-1} \mathbf{U} |$.

Then

$$\begin{aligned} | \mathbf{J} | &= \left| \frac{\partial(\mathbf{U}, \hat{\boldsymbol{\beta}}_G)}{\partial \mathbf{Y}} \right| = | \mathbf{B} \quad \mathbf{G} | = \left| \begin{bmatrix} \mathbf{B}^T \\ \mathbf{G}^T \end{bmatrix} [\mathbf{B} \quad \mathbf{G}] \right|^{1/2} \\ &= \left| \begin{bmatrix} \mathbf{B}^T \mathbf{B} & \mathbf{B}^T \mathbf{G} \\ \mathbf{G}^T \mathbf{B} & \mathbf{G}^T \mathbf{G} \end{bmatrix} \right|^{1/2} \\ &= | \mathbf{B}^T \mathbf{B} |^{1/2} | \mathbf{G}^T \mathbf{G} - \mathbf{G}^T \mathbf{B} (\mathbf{B}^T \mathbf{B})^{-1} \mathbf{B}^T \mathbf{G} |^{1/2} \\ &= 1 \times | \mathbf{G}^T \mathbf{G} - \mathbf{G}^T (\mathbf{I} - \mathbf{H}) \mathbf{G} |^{1/2} \\ &= | \mathbf{x}^T \mathbf{x} |^{-1/2} \neq 0 \end{aligned}$$

which implies that $(\mathbf{U}, \hat{\boldsymbol{\beta}}_G)$ is of full rank ($= N$). The vector $(\mathbf{U}, \hat{\boldsymbol{\beta}}_G)$ is a linear combination of normals and so is normal.

52

We have

$$p(\mathbf{y} \mid \boldsymbol{\alpha}, \boldsymbol{\beta}) = p(\mathbf{U}, \hat{\boldsymbol{\beta}}_G \mid \boldsymbol{\alpha}, \boldsymbol{\beta}) \mid \mathbf{J} \mid = p(\mathbf{U} \mid \hat{\boldsymbol{\beta}}_G, \boldsymbol{\alpha}, \boldsymbol{\beta}) p(\hat{\boldsymbol{\beta}}_G \mid \boldsymbol{\alpha}, \boldsymbol{\beta}) \mid \mathbf{J} \mid$$

and

$$\text{cov}(\mathbf{U}, \hat{\boldsymbol{\beta}}_G) = \text{E}[\mathbf{U}(\hat{\boldsymbol{\beta}}_G - \boldsymbol{\beta})^T] = \mathbf{0},$$

and so \mathbf{U} and $\hat{\boldsymbol{\beta}}_G$ are uncorrelated, and since normal therefore independent.

Hence

$$p(\mathbf{y} \mid \boldsymbol{\alpha}, \boldsymbol{\beta}) = p(\mathbf{U} \mid \boldsymbol{\alpha}) p(\hat{\boldsymbol{\beta}}_G \mid \boldsymbol{\alpha}, \boldsymbol{\beta}) \mid \mathbf{J} \mid.$$

Inference for $\boldsymbol{\lambda}$ may be based on the *marginal* likelihood

$$L_m(\boldsymbol{\lambda}) = p(\mathbf{s}_1 \mid \boldsymbol{\lambda}).$$

In the REML context we have $\mathbf{s}_1 = \mathbf{u}$, $\mathbf{s}_2 = \hat{\boldsymbol{\beta}}_G$, $\boldsymbol{\lambda} = \boldsymbol{\alpha}$, $\boldsymbol{\phi} = \boldsymbol{\beta}$, and $p(\mathbf{U} \mid \boldsymbol{\alpha})$ is a marginal likelihood.

Hence

$$p(\mathbf{U} \mid \boldsymbol{\alpha}) = \frac{p(\mathbf{y} \mid \boldsymbol{\alpha}, \boldsymbol{\beta})}{p(\hat{\boldsymbol{\beta}}_G \mid \boldsymbol{\alpha}, \boldsymbol{\beta})} \mid \mathbf{J} \mid^{-1}.$$

We have

$$p(\mathbf{y} \mid \boldsymbol{\alpha}, \boldsymbol{\beta}) = (2\pi)^{-N/2} \mid \mathbf{V} \mid^{-1/2} \exp \left\{ -\frac{1}{2} (\mathbf{y} - \mathbf{x}\boldsymbol{\beta})^T \mathbf{V}^{-1} (\mathbf{y} - \mathbf{x}\boldsymbol{\beta}) \right\},$$

53

and

$$\begin{aligned} p(\hat{\boldsymbol{\beta}}_G \mid \boldsymbol{\alpha}, \boldsymbol{\beta}) &= (2\pi)^{-(k+1)/2} \mid \mathbf{x}^T \mathbf{V}^{-1} \mathbf{x} \mid^{1/2} \\ &\times \exp \left\{ -\frac{1}{2} (\hat{\boldsymbol{\beta}}_G - \boldsymbol{\beta})^T \mathbf{x}^T \mathbf{V}^{-1} \mathbf{x} (\hat{\boldsymbol{\beta}}_G - \boldsymbol{\beta}) \right\} \end{aligned}$$

This leads to

$$\begin{aligned} p(\mathbf{U} \mid \boldsymbol{\alpha}) &= (2\pi)^{-(N-k-1)/2} \frac{\mid \mathbf{x}^T \mathbf{x} \mid^{1/2} \mid \mathbf{V} \mid^{-1/2}}{\mid \mathbf{x}^T \mathbf{V}^{-1} \mathbf{x} \mid^{1/2}} \\ &\times \exp \left\{ -\frac{1}{2} (\mathbf{y} - \mathbf{x}\hat{\boldsymbol{\beta}}_G)^T \mathbf{V}^{-1} (\mathbf{y} - \mathbf{x}\hat{\boldsymbol{\beta}}_G) \right\} \end{aligned} \quad (20)$$

which does not depend upon \mathbf{B} , hence we can choose any linear combination of the residuals.

- To summarize: the “data” \mathbf{U} (a linear combination of residuals from an OLS fit), has a distribution that depends on $\boldsymbol{\alpha}$ only – this defines a marginal likelihood (the REML likelihood) which may then be maximized as a function of $\boldsymbol{\alpha}$.
- The log marginal (restricted) likelihood is, upto a constant,

$$l_m(\boldsymbol{\alpha}) = -\frac{1}{2} \log |\mathbf{x}^T \mathbf{V}^{-1} \mathbf{x}| - \frac{1}{2} \log |\mathbf{V}| - \frac{1}{2} (\mathbf{y} - \mathbf{x} \hat{\boldsymbol{\beta}}_G)^T \mathbf{V}^{-1} (\mathbf{y} - \mathbf{x} \hat{\boldsymbol{\beta}}_G).$$

The profile log-likelihood based on \mathbf{Y} is:

$$l_P(\boldsymbol{\alpha}) = -\frac{1}{2} \log |\mathbf{V}| - \frac{1}{2} (\mathbf{y} - \mathbf{x} \hat{\boldsymbol{\beta}}_G)^T \mathbf{V}^{-1} (\mathbf{y} - \mathbf{x} \hat{\boldsymbol{\beta}}_G),$$

and so we have the additional term $-\frac{1}{2} \log |\mathbf{x}^T \mathbf{V} \mathbf{x}|$ that accounts for the degrees of freedom in estimation of $\boldsymbol{\beta}$.

- In terms of computation calculating REML estimators can be carried out with ML code, altered to include the extra term.

55

- In general, REML estimators have finite sample bias, but they are preferable to ML estimators, particularly for small samples.
- So far as estimation of the variance components are concerned, the asymptotic distribution of the ML/REML estimator is normal, with variance given by Fisher’s information.
- Suppose we fit two (nested) models using REML. Different sets of observations are used in each and so we cannot use a likelihood ratio on regression parameters to test whether the smaller model is a valid statistical simplification of the larger model.
- Likelihood ratio tests for variance components are valid.

56

Implementation of MLE and REML

MLE and REML require iteration between $\hat{\beta}|\hat{\alpha}$ and $\hat{\alpha}|\hat{\beta}$.

Originally the *EM algorithm* was used, e.g., Laird and Ware (1982, *Biometrics*). We illustrate for MLE and, for example, suppose $\mathbf{E}_i = \mathbf{I}_{n_i}\sigma^2$. The “missing data” here are the random effects \mathbf{b}_i and the errors ϵ_i .

The M-step: Given \mathbf{b}_i and ϵ_i , obtain estimates $\hat{\alpha} = (\hat{\sigma}^2, \hat{\mathbf{D}})$:

$$\begin{aligned}\hat{\sigma}^2 &= \frac{\sum_{i=1}^m \epsilon_i^T \epsilon_i}{\sum_{i=1}^m n_i} = \frac{t_1}{N} \\ \hat{\mathbf{D}} &= \frac{1}{m} \sum_{i=1}^m \mathbf{b}_i \mathbf{b}_i^T = \frac{\mathbf{t}_2}{m},\end{aligned}$$

where t_1 and \mathbf{t}_2 are the sufficient statistics.

The E step: Estimate the sufficient statistics given the current values $\hat{\alpha}$, via their expected values:

$$\begin{aligned}\hat{t}_1 &= \text{E} \left[\sum_{i=1}^m \epsilon_i^T \epsilon_i | \mathbf{y}_i, \hat{\beta}(\hat{\alpha}), \hat{\alpha} \right] \\ \hat{\mathbf{t}}_2 &= \text{E} \left[\sum_{i=1}^m \mathbf{b}_i^T \mathbf{b}_i | \mathbf{y}_i, \hat{\beta}(\hat{\alpha}), \hat{\alpha} \right].\end{aligned}$$

57

Closed form fixed and random effect estimates are available once we know α .

Slow convergence has been reported so that now the *Newton-Raphson method* is more frequently used.

Let θ be a $p \times 1$ parameter vector containing the variance components, $l(\cdot)$ the log-likelihood, \mathbf{G} the $p \times 1$ score vector, and $\mathbf{I}^*(\cdot)$ the $p \times p$ observed information matrix. Then a second order Taylor series expansion of $l(\cdot)$ about $\theta^{(t)}$, the estimate at iteration t gives:

$$\mathbf{g}^{(t)}(\theta) = l(\theta) + \mathbf{G}^{(t)T}(\theta - \theta^{(t)}) + \frac{1}{2}(\theta - \theta^{(t)})^T \mathbf{I}^{*(t)}(\theta - \theta^{(t)}),$$

differentiating and setting equal to zero:

$$\frac{\partial \mathbf{g}^{(t)}}{\partial \theta} = \mathbf{G}^{(t)} + \mathbf{I}^{*(t)}(\theta - \theta^{(t)}) = \mathbf{0},$$

gives the next estimate

$$\theta^{(t+1)} = \theta^{(t)} - \{\mathbf{I}^{*(t)}\}^{-1} \mathbf{G}^{(t)}.$$

The use of the expected information gives *Fisher's scoring method*.

See Lindstrom and Bates (1988, *JASA*) for details.

Lack of convergence of the algorithm/negative estimates, may sometimes indicate that a poor model is being fitted.

58

Dental Example

The simplest possible mixed effects model is given by

$$Y_{ij} = \beta_0 + b_i + \beta_1 t_j + \epsilon_{ij},$$

where ϵ_{ij} are iid with $E[\epsilon_{ij}] = 0$ and $\text{var}(\epsilon_{ij}) = \sigma_\epsilon^2$ and b_i represent random effects with $b_i \sim_{iid} N(0, \sigma_0^2)$, and represent perturbations for girl i from the population intercept β_0 .

Girl-specific intercepts $\beta_{0i} = \beta_0 + b_i$.

We could write b_{0i} , but use b_i for simplicity.

After conditioning on the random effect we have *independent* observations on each girl, we have assumed that allowing the intercepts to vary has removed all within-girl correlation.

59

The marginal distribution is normal with mean

$$E[\mathbf{Y}|\beta_0, \beta_1, \sigma_\epsilon^2, \sigma_0^2] = \boldsymbol{\mu},$$

where

$$\boldsymbol{\mu} = (\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_m)^T$$

is a $4m \times 1$ vector and

$$\boldsymbol{\mu}_1 = \dots = \boldsymbol{\mu}_m = (\beta_0 + \beta_1 t_1, \beta_0 + \beta_1 t_2, \beta_0 + \beta_1 t_3, \beta_0 + \beta_1 t_4)^T.$$

The variance is given by

$$\text{var}(\mathbf{Y}|\beta_0, \beta_1, \sigma_\epsilon^2, \sigma_0^2) = \mathbf{V},$$

where \mathbf{V} is the $4m \times 4m$ block diagonal matrix with

$$\mathbf{V}_i = \text{var}(\mathbf{Y}_i) = \sigma^2[\mathbf{J}_{n_i}\rho + \mathbf{I}_{n_i}(1 - \rho)],$$

with $\sigma^2 = \sigma_\epsilon^2 + \sigma_0^2$ and $\rho = \frac{\sigma_0^2}{\sigma^2} = \frac{\sigma_0^2}{\sigma_\epsilon^2 + \sigma_0^2}$. Hence the random intercepts model induces a marginal form with constant variances and constant correlations on measurements on the same child, regardless of the time between observations.

We analyze the dental data using LMEMs. To do this we use the `nlme` package which is described in Pinheiro and Bates (2000) – very flexible, but the syntax is not always obvious...

The `groupedData` function is useful for plotting and modeling (attaches a model function as an attribute to a dataset).

```
> library(nlme)
> data(Orthodont) # Dental data is one of the data sets in the package.
> Orthgirl <- Orthodont[Orthodont$Sex=="Female",]
> trellldat <- groupedData( distance ~ age | Subject, data=Orthgirl )
> plot(trellldat)
```

Figure 3 shows the data plotted using a “trellis” plot – note that data are not plotted in the original order.

61

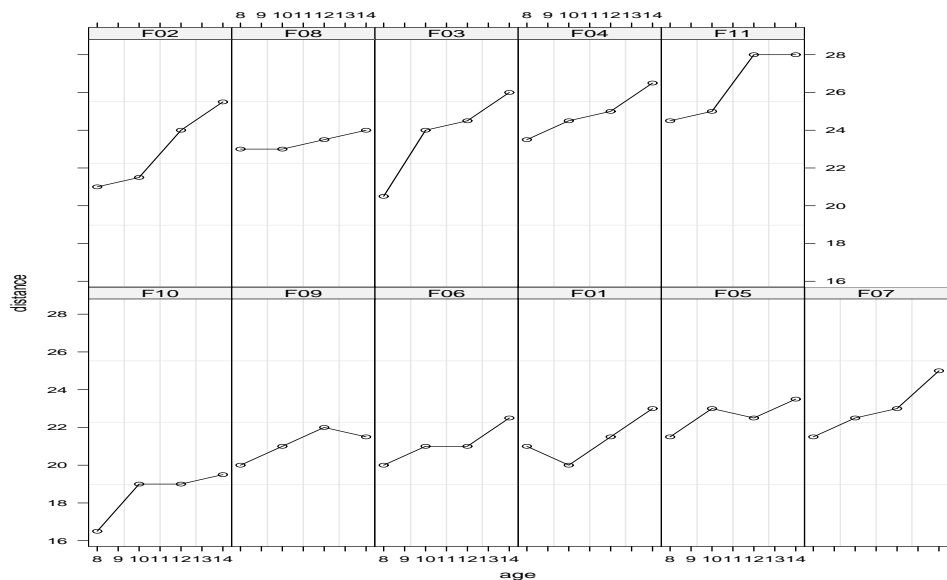


Figure 3: Length versus age (in years) for 11 girls.

62

We now carry out parameter estimation, first naively, and then using LMEM via REML.

```
> summary(lm(distance~age,data=Orthgirl))
```

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	17.3727	1.6378	10.608	1.87e-13 ***
age	0.4795	0.1459	3.287	0.00205 **

```
> summary(lme( distance ~ age, data = Orthgirl, random = ~1 | Subject ))
```

Linear mixed-effects model fit by REML

Random effects:

Formula: ~1 | Subject

	(Intercept)	Residual
StdDev:	2.06847	0.7800331

Fixed effects: distance ~ age

	Value	Std.Error	DF	t-value	p-value
(Intercept)	17.372727	0.8587419	32	20.230440	0
age	0.479545	0.0525898	32	9.118598	0

63

Notice the standard error for β_1 is smaller for the REML analysis – slopes are being estimated from within-girl comparisons.

The REML estimates of the variance components are $\hat{\sigma}_\epsilon = 0.78$, $\hat{\sigma}_0 = 2.07$ so that $\hat{\rho} = 0.875$ which ties in with the empirical correlations (16). The marginal standard deviation is given by $(\hat{\sigma}_\epsilon^2 + \hat{\sigma}_0^2)^{1/2} = 2.21$, in agreement with the diagonal elements of (16).

64

Now for comparison we fit the LMEM with ML:

```
> summary(lme( distance ~ age, data = Orthgirl, random = ~1 | Subject, method = "ML" )
Linear mixed-effects model fit by maximum likelihood
Random effects:
  Formula: ~1 | Subject
          (Intercept)  Residual
StdDev:      1.969870  0.7681235
Fixed effects: distance ~ age
              Value Std.Error DF   t-value p-value
(Intercept) 17.372727 0.8506287 32 20.423397      0
age          0.479545 0.0530056 32  9.047078      0
```

Note that the MLEs of the variance components are smaller than the REML counterparts. Slight differences in the standard errors of the fixed effects (but not a big difference here).

65

Bayesian Justification for REML

Another justification is to assign a flat improper prior to the regression coefficients and then integrate these from the model.

Example: Normal Linear Model

Consider the linear regression for independent data: $\mathbf{Y}|\boldsymbol{\beta}, \sigma^2 \sim N(\mathbf{x}\boldsymbol{\beta}, \mathbf{I}_n\sigma^2)$, with $\dim(\boldsymbol{\beta}) = k + 1$.

Consider

$$p(\mathbf{y}|\sigma^2) = \int p(\mathbf{y}|\boldsymbol{\beta}, \sigma^2)\pi(\boldsymbol{\beta})d\boldsymbol{\beta},$$

and assume $\pi(\boldsymbol{\beta}) \propto 1$.

66

Hence

$$\begin{aligned}
 p(\mathbf{y}|\sigma^2) &= \int (2\pi\sigma^2)^{-n/2} \exp \left[-\frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{x}\boldsymbol{\beta})^\top (\mathbf{y} - \mathbf{x}\boldsymbol{\beta}) \right] d\boldsymbol{\beta} \\
 &= (2\pi\sigma^2)^{-n/2} \int \exp \left[-\frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{x}\hat{\boldsymbol{\beta}} + \mathbf{x}\hat{\boldsymbol{\beta}} - \mathbf{x}\boldsymbol{\beta})^\top \right. \\
 &\quad \times \left. (\mathbf{y} - \mathbf{x}\hat{\boldsymbol{\beta}} + \mathbf{x}\hat{\boldsymbol{\beta}} + \mathbf{x}\boldsymbol{\beta}) \right] d\boldsymbol{\beta} \\
 &= (2\pi\sigma^2)^{-(n-k-1)/2} \exp \left[-\frac{RSS}{2\sigma^2} \right] |\mathbf{x}^\top \mathbf{x}|^{-1/2}
 \end{aligned}$$

where the residual sum of squares

$$RSS = (\mathbf{y} - \mathbf{x}\hat{\boldsymbol{\beta}})^\top (\mathbf{y} - \mathbf{x}\hat{\boldsymbol{\beta}}).$$

Maximization of $l(\sigma^2) = \log p(\mathbf{y}|\sigma^2)$ yields the unbiased estimator

$$\hat{\sigma}^2 = \frac{RSS}{n - k - 1}.$$

67

Example: LMEM

Again obtain the distribution of the data as a function of $\boldsymbol{\alpha}$ only, by integrating $\boldsymbol{\beta}$ from the model, and assuming an improper flat prior for $\boldsymbol{\beta}$.

We have

$$p(\mathbf{y}|\boldsymbol{\alpha}) = \int_{\boldsymbol{\beta}} p(\mathbf{y}|\boldsymbol{\beta}, \boldsymbol{\alpha}) \times \pi(\boldsymbol{\beta}) d\boldsymbol{\beta},$$

leading to

$$\begin{aligned}
 l(\boldsymbol{\alpha}) &= \log p(\mathbf{y}|\boldsymbol{\alpha}) = -\frac{1}{2} \sum_{i=1}^m \log |\mathbf{V}_i(\boldsymbol{\alpha})| \\
 &\quad - \frac{1}{2} \sum_{i=1}^m \log |\mathbf{x}_i^\top \mathbf{V}_i(\boldsymbol{\alpha}) \mathbf{x}_i| - \frac{1}{2} \sum_{i=1}^m (\mathbf{y}_i - \mathbf{x}_i \hat{\boldsymbol{\beta}})^\top \mathbf{V}_i^{-1}(\boldsymbol{\alpha}) (\mathbf{y}_i - \mathbf{x}_i \hat{\boldsymbol{\beta}}),
 \end{aligned}$$

which differs from the “usual” likelihood by the term

$$-\frac{1}{2} \sum_{i=1}^m \log |\mathbf{x}_i^\top \mathbf{V}_i^{-1}(\boldsymbol{\alpha}) \mathbf{x}_i|.$$

This expression is the same as that which results from the maximization of the distribution of the residuals.

Estimates of $\boldsymbol{\beta}$ change since they are a function of $\hat{\boldsymbol{\alpha}}$.

68

Inference for Random Effects

Examples:

- Pharmacokinetics: individualization of a profile.
- Dairy herds: genetic merit of a particular bull – data are in the form of the milk yields of his daughters.
- Psychology: inference for the IQ of an individual from a set of test scores.
- Industrial applications: operating characteristics of a particular machine.

From a frequentist perspective, inference for random effects is often viewed as *prediction* rather than estimation, since \mathbf{b} are random variables.

The usual frequentist optimality criteria for a fixed effect $\boldsymbol{\theta}$, are based upon unbiasedness:

$$E[\hat{\boldsymbol{\theta}}] - \boldsymbol{\theta} = \mathbf{0},$$

where $\boldsymbol{\theta}$ is a fixed constant, and upon the variance of the estimator

$$\text{var}(\hat{\boldsymbol{\theta}}).$$

These need to be adjusted when inference is required for a random effect \mathbf{b} .

69

We wish to find a predictor $\tilde{\mathbf{b}} = f(\mathbf{Y})$ of \mathbf{b} .

An unbiased predictor $\tilde{\mathbf{b}}$ is such that

$$E_{y,b}[\tilde{\mathbf{b}} - \mathbf{b}] = E[\tilde{\mathbf{b}} - \mathbf{b}] = \mathbf{0},$$

to give

$$E[\tilde{\mathbf{b}}] = E[\mathbf{b}]$$

so that the expectation of the predictor is equal to the expectation of the random variable that it is predicting.

The variance of a random variable is defined with respect to a fixed number, the mean. In the context of prediction of a random variability, a more relevant summary of the variability is

$$\text{var}(\tilde{\mathbf{b}} - \mathbf{b}) = \text{var}(\tilde{\mathbf{b}}) + \text{var}(\mathbf{b}) - 2\text{cov}(\tilde{\mathbf{b}}, \mathbf{b}).$$

There are many different criteria that may be used to find a predictor.

Since we are predicting a random variable it is natural to use minimum mean squared error (MSE) as a criteria, rather than minimum variance.

The MSE of $\tilde{\mathbf{b}}$ is given by

$$\text{MSE}(\tilde{\mathbf{b}}) = E_{y,b}[(\tilde{\mathbf{b}} - \mathbf{b})^T \mathbf{A}(\tilde{\mathbf{b}} - \mathbf{b})],$$

for non-singular \mathbf{A} .

This leads to $\tilde{\mathbf{b}} = E[\mathbf{b} | \mathbf{y}]$, irrespective of \mathbf{A} (see Exercises 2). Hence the best prediction is that which estimates the random variable by its conditional mean.

We now examine properties of $\tilde{\mathbf{b}}$.

Unbiasedness

We have

$$E_y[\tilde{\mathbf{b}}] = E_y\{E_{b|y}[\mathbf{b} | \mathbf{y}]\} = E_b[\mathbf{b}]$$

where we first step follows on substitution of $\tilde{\mathbf{b}}$ and the second from iterated expectation. (Note: $E_u[U] = E_{u,v}[U] = E_v\{E_{u|v}[U|V]\}$.)

71

Variability

Recall an appropriate measure of variability:

$$\text{var}(\tilde{\mathbf{b}}_i - \mathbf{b}_i) = \text{var}(\tilde{\mathbf{b}}_i) + \text{var}(\mathbf{b}_i) - 2\text{cov}(\tilde{\mathbf{b}}_i, \mathbf{b}_i).$$

We have

$$\begin{aligned} \text{cov}_{\tilde{\mathbf{b}}, \mathbf{b}}(\tilde{\mathbf{b}}_i, \mathbf{b}_i) &= E_y[\text{cov}(\tilde{\mathbf{b}}_i, \mathbf{b}_i | \mathbf{y})] + \text{cov}_y(E[\tilde{\mathbf{b}}_i | \mathbf{y}], E[\mathbf{b}_i | \mathbf{y}]) \\ &= E_y[\text{cov}(\tilde{\mathbf{b}}_i, \mathbf{b}_i | \mathbf{y})] + \text{cov}_y(\tilde{\mathbf{b}}_i, \tilde{\mathbf{b}}_i) \\ &= \text{var}(\tilde{\mathbf{b}}_i) \end{aligned} \tag{21}$$

The first term in (21) is the covariance between a constant $E[\tilde{\mathbf{b}} | \mathbf{y}]$ (since \mathbf{y} is conditioned upon), and $\tilde{\mathbf{b}}$, and so is zero (because the covariance between a constant and any quantity is zero). In the second term we have used $E[\tilde{\mathbf{b}}_i | \mathbf{y}] = E[E[\mathbf{b}_i | \mathbf{y}] | \mathbf{y}] = \tilde{\mathbf{b}}_i$.

Hence

$$\text{var}(\tilde{\mathbf{b}}_i - \mathbf{b}_i) = \text{var}(\mathbf{b}_i) - \text{var}(\tilde{\mathbf{b}}_i).$$

72

Application to the LMEM

The predictor, $\tilde{\mathbf{b}} = E[\mathbf{b} | \mathbf{y}]$, is a random variable, since it is a function of \mathbf{y} , and so we need to know something about $p(\mathbf{b} | \mathbf{y})$ in order to derive its form.

Definitions: Suppose \mathbf{U} is an $n \times 1$ vector of random variables, and \mathbf{V} is an $m \times 1$ vector of random variables. Then $\text{cov}(\mathbf{U}, \mathbf{V}) = \mathbf{C}$ is an $n \times m$ matrix with (i, j) -th element $\text{cov}(U_i, V_j)$, $i = 1, \dots, n; j = 1, \dots, m$. Also $\text{cov}(\mathbf{V}, \mathbf{U}) = \mathbf{C}^T$. Now suppose $\mathbf{V} = \mathbf{AU}$ where \mathbf{A} is an $m \times n$ matrix. Then $\text{cov}(\mathbf{U}, \mathbf{AU}) = \mathbf{WA}^T$ where $\mathbf{W} = \text{cov}(\mathbf{U})$, and $\text{cov}(\mathbf{AU}, \mathbf{U}) = \mathbf{AW}$.

Consider the LMEM

$$\mathbf{y} = \mathbf{x}\boldsymbol{\beta} + \mathbf{z}\mathbf{b} + \boldsymbol{\epsilon},$$

and assume \mathbf{b} and $\boldsymbol{\epsilon}$ are independent and $\mathbf{b}_i \sim N(\mathbf{0}, \mathbf{D})$, $\boldsymbol{\epsilon} \sim N(\mathbf{0}, \sigma_\epsilon^2 \mathbf{I})$ then, using the above results:

$$\begin{bmatrix} \mathbf{b}_i \\ \mathbf{y}_i \end{bmatrix} \sim N_{q+1+n_i} \left(\begin{bmatrix} \mathbf{0} \\ \mathbf{x}_i \boldsymbol{\beta} \end{bmatrix}, \begin{bmatrix} \mathbf{D} & \mathbf{D}\mathbf{z}_i^T \\ \mathbf{z}_i \mathbf{D} & \mathbf{V}_i \end{bmatrix} \right).$$

since

$$\text{cov}(\mathbf{b}_i, \mathbf{y}_i) = \text{cov}(\mathbf{b}_i, \mathbf{x}_i \boldsymbol{\beta} + \mathbf{z}_i \mathbf{b}_i + \boldsymbol{\epsilon}_i) = \text{cov}(\mathbf{b}_i, \mathbf{z}_i \mathbf{b}_i) = \mathbf{D}\mathbf{z}_i^T,$$

and similarly $\text{cov}(\mathbf{y}_i, \mathbf{b}_i) = \mathbf{z}_i \mathbf{D}$.

73

Using properties of the multivariate normal distribution, the predictor takes the form:

$$\tilde{\mathbf{b}}_i = E[\mathbf{b}_i | \mathbf{y}_i] = \mathbf{D}\mathbf{z}_i^T \mathbf{V}_i^{-1} (\mathbf{y}_i - \mathbf{x}_i \boldsymbol{\beta}) \quad (22)$$

This is known as the best linear unbiased predictor (BLUP), where unbiased refers to it satisfying $E[\tilde{\mathbf{b}}_i] = E[\mathbf{b}_i]$.

The random effect predictor is a shrinkage estimator since it pulls the fixed effect estimator towards zero, as we see in examples later.

The form (22) is not of practical use since it depends on the unknown $\boldsymbol{\beta}$ and $\boldsymbol{\alpha}$; instead we use

$$\tilde{\mathbf{b}}_i = E[\mathbf{b}_i | \mathbf{y}_i] = \hat{\mathbf{D}}\mathbf{z}_i^T \hat{\mathbf{V}}_i^{-1} (\mathbf{y}_i - \mathbf{x}_i \hat{\boldsymbol{\beta}}). \quad (23)$$

Substitution of $\hat{\boldsymbol{\beta}}$ is not such a problem (since it is an unbiased estimator, and appears in (22) in a linear fashion), but $\hat{\boldsymbol{\alpha}}$ is more problematic.

The uncertainty in the prediction is given by

$$\text{var}(\tilde{\mathbf{b}}_i - \mathbf{b}_i) = \text{var}(\mathbf{b}_i) - \text{var}(\tilde{\mathbf{b}}_i) = \mathbf{D} - \text{var}(\tilde{\mathbf{b}}_i)$$

We have

$$\tilde{\mathbf{b}}_i = \mathbf{D}\mathbf{z}_i^T \mathbf{V}_i^{-1} (\mathbf{Y}_i - \mathbf{x}_i \hat{\boldsymbol{\beta}}) = \mathbf{K}_i (\mathbf{Y}_i - \mathbf{x}_i \hat{\boldsymbol{\beta}}),$$

and

$$\text{var}(\mathbf{Y}_i - \mathbf{x}_i \hat{\boldsymbol{\beta}}) = \text{var}(\mathbf{Y}_i) + \mathbf{x}_i \text{var}(\hat{\boldsymbol{\beta}}) \mathbf{x}_i^T - 2\text{cov}(\mathbf{Y}_i, \mathbf{x}_i \hat{\boldsymbol{\beta}}).$$

Since

$$\hat{\boldsymbol{\beta}} = (\mathbf{x}^T \mathbf{V}^{-1} \mathbf{x})^{-1} \sum_{i=1}^m \mathbf{x}_i^T \mathbf{V}_i^{-1} \mathbf{Y}_i,$$

we have

$$\text{cov}(\mathbf{Y}_i, \mathbf{x}_i \hat{\boldsymbol{\beta}}) = \mathbf{x}_i (\mathbf{x}^T \mathbf{V}^{-1} \mathbf{x})^{-1} \mathbf{x}_i^T \mathbf{V}_i^{-1} \text{var}(\mathbf{Y}_i) = \mathbf{x}_i \text{var}(\hat{\boldsymbol{\beta}}) \mathbf{x}_i^T,$$

and so

$$\text{var}(\tilde{\mathbf{b}}_i) = \mathbf{K}_i [\text{var}(\mathbf{Y}_i) - \mathbf{x}_i \text{var}(\hat{\boldsymbol{\beta}}) \mathbf{x}_i^T] \mathbf{K}_i^T = \mathbf{K}_i [\mathbf{V}_i - \mathbf{x}_i \text{var}(\hat{\boldsymbol{\beta}}) \mathbf{x}_i^T] \mathbf{K}_i^T$$

to give

$$\text{var}(\tilde{\mathbf{b}}_i - \mathbf{b}_i) = \mathbf{D} - \mathbf{D}\mathbf{z}_i^T \mathbf{V}_i^{-1} \mathbf{z}_i \mathbf{D} + \mathbf{D}\mathbf{z}_i^T \mathbf{V}_i^{-1} \mathbf{x}_i (\mathbf{x}^T \mathbf{V}^{-1} \mathbf{x})^{-1} \mathbf{x}_i^T \mathbf{V}_i^{-1} \mathbf{z}_i \mathbf{D}.$$

The variability of the prediction does not acknowledge the uncertainty in $\hat{\boldsymbol{\alpha}}$.

75

We now examine fitted values:

$$\begin{aligned} \hat{\mathbf{Y}}_i &= \mathbf{x}_i \hat{\boldsymbol{\beta}} + \mathbf{z}_i \hat{\mathbf{b}}_i \\ &= \mathbf{x}_i \hat{\boldsymbol{\beta}} + \mathbf{z}_i \{ \mathbf{D}\mathbf{z}_i^T \mathbf{V}_i^{-1} (\mathbf{Y}_i - \mathbf{x}_i \hat{\boldsymbol{\beta}}) \} \\ &= (\mathbf{I}_{n_i} - \mathbf{z}_i \mathbf{D}\mathbf{z}_i^T \mathbf{V}_i^{-1}) \mathbf{x}_i \hat{\boldsymbol{\beta}} + \mathbf{z}_i \mathbf{D}\mathbf{z}_i^T \mathbf{V}_i^{-1} \mathbf{Y}_i, \end{aligned}$$

a weighted combination of the population profile, and the unit's data.

Note that if $\mathbf{D} = \mathbf{0}$ we obtain $\hat{\mathbf{Y}}_i = \mathbf{x}_i \hat{\boldsymbol{\beta}}$.

We can also write

$$\hat{\mathbf{Y}}_i = \sigma_\epsilon^2 \mathbf{V}_i^{-1} \mathbf{x}_i \hat{\boldsymbol{\beta}} + (\mathbf{I}_{n_i} - \sigma_\epsilon^2 \mathbf{V}_i^{-1}) \mathbf{Y}_i$$

so that as $\sigma_\epsilon^2 \rightarrow 0$, $\hat{\mathbf{Y}}_i \rightarrow \mathbf{Y}_i$.

76

Example: One-way ANOVA

For the simple balanced ANOVA model previously considered

$$\tilde{b}_i = \frac{n\sigma_0^2}{\sigma_\epsilon^2 + n\sigma_0^2}(\bar{y}_i - \beta_0).$$

In practice we have an estimate $\hat{\beta}_0$, and the predictor is a weighted combination of the distance $\bar{y}_i - \hat{\beta}_0$ and zero. Hence for finite n the predictor is biased towards zero (recall our definition of unbiasedness is in terms of \mathbf{b}).

As $n \rightarrow \infty$, $\tilde{b}_i \rightarrow \bar{y}_i - \hat{\beta}_0$, so that

$$\hat{\beta}_0 + \tilde{b}_i \rightarrow \bar{y}_i \rightarrow E[Y_i].$$

77

The form

$$\tilde{\mathbf{b}}_i = E[\mathbf{b}_i | \mathbf{y}_i] = \hat{\mathbf{D}} \mathbf{z}_i^T \hat{\mathbf{V}}_i^{-1} (\mathbf{y}_i - \mathbf{x}_i \hat{\boldsymbol{\beta}})$$

can be justified in a number of ways, other than MSE.

Rather than assume normality we could consider estimators that are *linear* in \mathbf{y} . In Exercises 2 we show that this again leads to the above form.

Hence the best linear predictor is identical to the best predictor under normality.

For general distributions, $E[\mathbf{b}_i | \mathbf{y}_i]$ is not necessarily linear in \mathbf{y} . Once we plug $\boldsymbol{\alpha}$ into the BLUP we don't even have a linear predictor.

The BLUP is an empirical Bayes estimator. We should be considering $E[\mathbf{b} | \mathbf{y}]$, with

$$p(\mathbf{b} | \mathbf{y}) = \int \int p(\mathbf{b}, \boldsymbol{\beta}, \boldsymbol{\alpha} | \mathbf{y}) d\boldsymbol{\beta} d\boldsymbol{\alpha} = \int \int p(\mathbf{b} | \boldsymbol{\beta}, \boldsymbol{\alpha}, \mathbf{y}) p(\boldsymbol{\beta}, \boldsymbol{\alpha} | \mathbf{y}) d\boldsymbol{\beta} d\boldsymbol{\alpha},$$

but instead the BLUP is the mean of the distribution

$$p(\mathbf{b} | \hat{\boldsymbol{\beta}}, \hat{\boldsymbol{\alpha}}, \mathbf{y}),$$

so that rather than integrating over $\boldsymbol{\beta}, \boldsymbol{\alpha}$, estimates have been conditioned upon.

78

Example: Dental Growth

We again fit a LMEM with random intercepts only.

```
> remlelm <- lme(distance~I(age-11),data = Orthgirl,random = ~1 | Subject)
> summary(remlelm)
Formula: ~1 | Subject
          (Intercept)  Residual
StdDev:      2.06847  0.7800331
              Value Std.Error DF t-value p-value
(Intercept) 22.647727 0.6346568 32 35.6850      0
I(age - 11)  0.479545 0.0525898 32  9.1186      0
> b0hat <- b1hat <- NULL
> for (i in 1:11){
  x <- Orthgirl$age[seq((i-1)*4+1,(i-1)*4+4)]-11
  y <- Orthgirl$distance[seq((i-1)*4+1,(i-1)*4+4)]
  mod <- lm(y~x)
  b0hat[i] <- mod$coef[1]
  b1hat[i] <- mod$coef[2]
}
> index <- c(10,9,6,1,5,7,2,8,3,4,11)
> LSb0hat <- b0hat[index]; LSb1hat <- b1hat[index]
```

79

Shrinkage of Intercepts

```
> cbind(LSb0hat,LSb1hat,rcoef)
      LSb0hat LSb1hat (Intercept) I(age - 11)
F10  18.500   0.450   18.64240   0.4795455
F09  21.125   0.275   21.17728   0.4795455
F06  21.125   0.375   21.17728   0.4795455
F01  21.375   0.375   21.41869   0.4795455
F05  22.625   0.275   22.62578   0.4795455
F07  23.000   0.550   22.98791   0.4795455
F02  23.000   0.800   22.98791   0.4795455
F08  23.375   0.175   23.35003   0.4795455
F03  23.750   0.850   23.71216   0.4795455
F04  24.875   0.475   24.79853   0.4795455
F11  26.375   0.675   26.24704   0.4795455
```

Note ordering difference in coefficients from `lme`, and the slight shrinkage here towards the overall mean of 22.65; not much shrinkage here since $\hat{\sigma}_0$ is large relative to $\hat{\sigma}_\epsilon$ (see Figure 4).

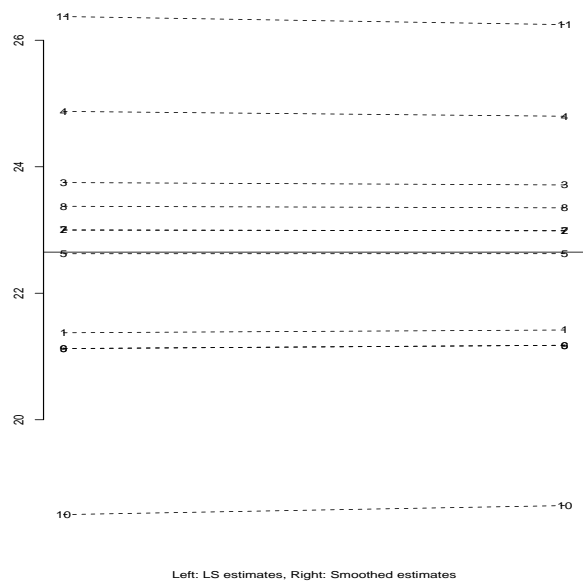


Figure 4: Least squares estimates and smoothed estimates, $\hat{\beta}_0 + \tilde{b}_i$.

Dental Example: Boys and Girls Joint Analyses

Table 2 describes LMEMs applied to the dental data and Table 3 results.

Model	Description											
1	Separate fits, random intercepts											
2	Separate fits, random intercepts and slopes, uncorrelated											
3	Separate fits, random intercepts and slopes, correlated											
4	Combined fit, separate intercepts, common slope, random intercepts											
5	Combined fit, separate intercepts and slopes, random intercepts											
6	Combined fit, separate intercepts and slopes, random intercepts and slopes, uncorrelated											
7	Combined fit, separate intercepts and slopes, random intercepts and slopes, correlated											

Table 2: Various LMEMs.

Model	Boys						Girls					
	$\hat{\beta}_0$	$\hat{\beta}_1$	$\hat{\sigma}_0$	$\hat{\sigma}_1$	$\hat{\rho}_{01}$	$\hat{\sigma}_\epsilon$	$\hat{\beta}_0$	$\hat{\beta}_1$	$\hat{\sigma}_0$	$\hat{\sigma}_1$	$\hat{\rho}_{01}$	$\hat{\sigma}_\epsilon$
1	25.0	0.78	1.63	—	—	1.68	22.7	0.48	2.07	—	—	0.78
2	25.0	0.78	1.64	0.19	—	1.61	22.6	0.48	2.08	0.16	—	0.67
3	25.0	0.78	1.64	0.19	-0.01	1.61	22.6	0.48	2.08	0.16	0.53	0.67
4	25.0	0.66	1.81	—	—	1.43	22.6	0.66	1.81	—	—	1.43
5	25.0	0.78	1.82	—	—	1.39	22.6	0.48	1.82	—	—	1.39
6	25.0	0.78	1.83	0.18	—	1.31	22.6	0.48	1.83	0.18	—	1.31
7	25.0	0.78	1.83	0.18	0.21	1.31	22.6	0.48	1.83	0.18	0.21	1.31

Table 3: Various LMEM analyses.

R code for models

```
# Set parameterization (to corner point)
> options(contrasts=c("contr.treatment","contr.poly"))
# Separate fits - intercept only, model 1
> remlF <- lme( distance ~ I(age-11), data = Orthgirl, random = ~1 )
> remlM <- lme( distance ~ I(age-11), data = Orthboy, random = ~1 )
# Separate fits - intercept and age, diagonal, model 2
> remlF2d <- lme( distance ~ I(age-11), data = Orthgirl, random = pdDiag(~I(age-11)))
> remlM2d <- lme( distance ~ I(age-11), data = Orthboy, random = pdDiag(~I(age-11)))
# Separate fits - intercept and age, non-diagonal, model 3
> remlF2 <- lme( distance ~ I(age-11), data = Orthgirl, random = ~I(age-11))
> remlM2 <- lme( distance ~ I(age-11), data = Orthboy, random = ~I(age-11))
# Combined fit - common slope, intercept only, model 4
> remlMF <- lme( distance ~ I(age-11)+Sex, data = Orthodont, random = ~1 )
# Combined fit - separate intercepts and slopes, intercept only - model 5
> remlMFi <- lme( distance ~ I(age-11)+Sex+I(age-11):Sex, data = Orthodont,
  random = ~1 )
# Combined fit - sep intercepts and slopes, uncor random intercepts and slopes - model 6
> remlMF2 <- lme( distance ~ I(age-11)+Sex+I(age-11):Sex, data = Orthodont,
  random=pdDiag(~I(age-11)) )
# Combined fit - sep intercepts and slopes, cor random intercepts and slopes - model 7
> remlMF3 <- lme( distance ~ I(age-11)+Sex+I(age-11):Sex, data = Orthodont,
  random=~I(age-11) )
```

83

Example of Output (model 4)

```
> summary(remlMF)
Random effects:
Formula: ~1 | Subject
(Intercept) Residual
StdDev:    1.807425 1.431592
Fixed effects: distance ~ I(age - 11) + Sex
              Value Std.Error DF   t-value p-value
(Intercept) 24.968750 0.4860008 80 51.37595 0.0000
I(age - 11)  0.660185 0.0616059 80 10.71626 0.0000
SexFemale   -2.321023 0.7614168 25 -3.04829 0.0054
Correlation:
      (Intr) I(-11)
I(age - 11)  0.000
SexFemale   -0.638 0.000
Number of Observations: 108
Number of Groups: 27
```

Figure 5 gives normal QQ plots of the LS estimates of intercepts and slopes, for boys and girls.

Figure 6 gives a scatter plot of the LS estimates of intercepts and slopes, for boys and girls.

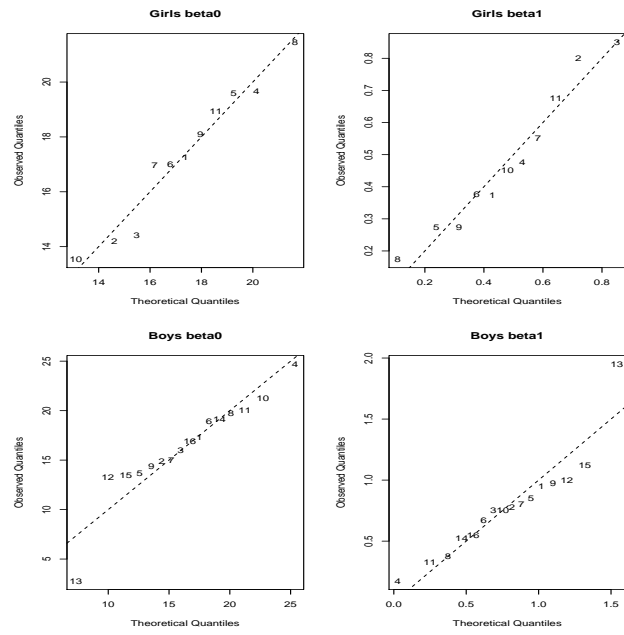


Figure 5: QQ plot of the LS estimates.

85

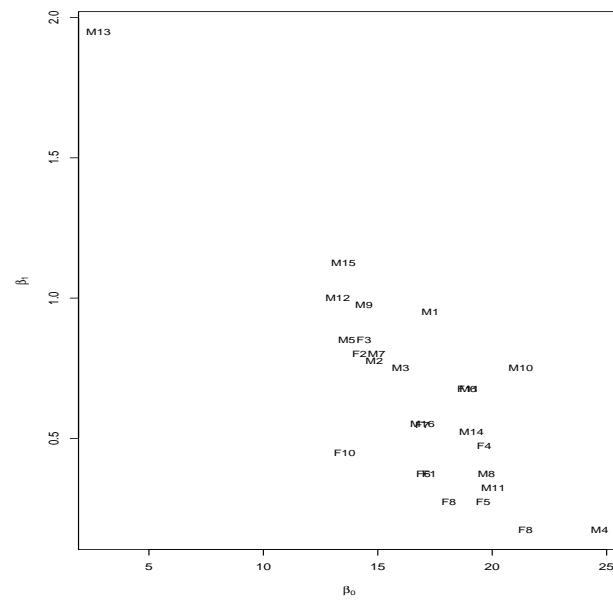


Figure 6: Plot of the LS estimates for boys and girls.

86