

## CHAPTER 8: LINEAR MODELS

The effect of ignoring dependence is two-fold. First, standard errors reported from independent data methods are calculated under the assumption of independence and so, if the data are truly dependent, will be inaccurate. Second, and more subtly, models for dependence may control for confounding, e.g. by time in the air pollution example.

While making inference for dependent data is more difficult than for independent data, designs that collect dependent data can be very efficient. For example, 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).

For longitudinal data there are two extreme fixed effects approaches. Proceeding naively we could assume a single “marginal” curve for *all* of the data, and carry out a standard analysis assuming independent data.

21

### *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

$$\begin{bmatrix} \sigma_1 & & & & \\ \rho_{12} & \sigma_2 & & & \\ \rho_{13} & \rho_{23} & \sigma_3 & & \\ \rho_{14} & \rho_{24} & \rho_{34} & \sigma_4 & \end{bmatrix} \quad (9)$$

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,  $j \neq k, j, k = 1, \dots, 4$ .

22

Across girls we may empirically estimate the entries of (9) by

$$\begin{bmatrix} 2.12 \\ 0.83 & 1.90 \\ 0.86 & 0.90 & 2.36 \\ 0.84 & 0.88 & 0.95 & 2.44 \end{bmatrix} \quad (10)$$

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

23

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

While providing valid inference for each curve, there is no "borrowing of strength" across children, that is, each girl's fit is based only on their data alone, and not on those of other girls.

We would hope that if there is *similarity* between the curves, and that the totality of data will aid in the estimation of each individual curve. In some instances this may be vital, for example, if  $n_i = 1$  for a particular individual, then their own data alone will not allow parameter estimation.

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

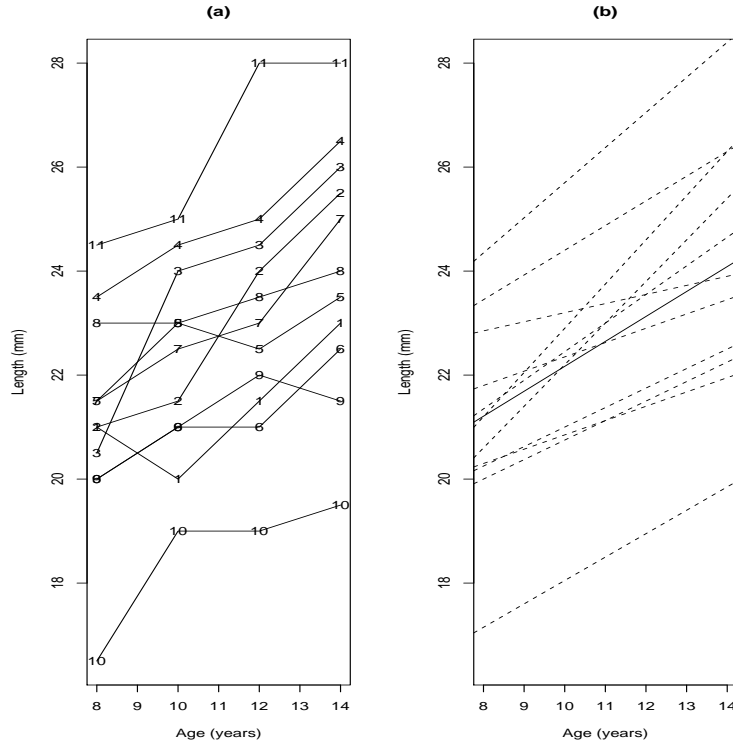


Figure 5: Dental plots for girls only: (a) Individual observed data (with plotting symbol girl index), (b) Individual fitted curves (dashed) and overall fitted curve

25

### 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  $\mathbf{z}_{ij} = (1, z_{ij1}, \dots, z_{ijq})^T$ .

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

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

$$\mathbf{y}_i = \mathbf{x}_i \boldsymbol{\beta} + \mathbf{z}_i \mathbf{b}_i + \boldsymbol{\epsilon}_i, \quad (11)$$

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

*Stage 2:* Model for random terms:

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

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

The two stages define the marginal model:

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

We describe likelihood and Bayesian approaches to inference.

27

### 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

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

Here  $\boldsymbol{\alpha} = (\sigma_\epsilon^2, \mathbf{D})$  denote the variance-covariance parameters. Here  $\mathbf{V} = \mathbf{z} \mathbf{D} \mathbf{z}^T + \sigma_\epsilon^2 \mathbf{I}_N$ , 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:

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

28

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

$$\begin{aligned} l(\boldsymbol{\beta}, \boldsymbol{\alpha}) = & - \frac{N}{2} \log 2\pi - \frac{1}{2} \sum_{i=1}^m \log |\mathbf{V}_i(\boldsymbol{\alpha})| \\ & - \frac{1}{2} \sum_{i=1}^m (\mathbf{Y}_i - \mathbf{x}_i\boldsymbol{\beta})^\top \mathbf{V}_i(\boldsymbol{\alpha})^{-1} (\mathbf{Y}_i - \mathbf{x}_i\boldsymbol{\beta}). \end{aligned} \quad (12)$$

29

### Example: One-way ANOVA

Consider the simple ANOVA model

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

with  $b_i$  and  $\epsilon_{ij}$  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  $E[\mathbf{Y}] = \mathbf{1}_N \beta_0$ ,  $\text{var}(\mathbf{Y}) = \mathbf{V} = \mathbf{1}_N \mathbf{1}_N^\top \sigma_0^2 + \mathbf{I}_N \sigma_\epsilon^2 = \mathbf{J}_N \sigma_0^2 + \mathbf{I}_N \sigma_\epsilon^2$ .

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

$$\sigma^2 \begin{bmatrix} 1 & \rho & \rho & \rho & \cdot & \cdot & \cdot & \cdot & 0 & 0 & 0 & 0 \\ \rho & 1 & \rho & \rho & \cdot & \cdot & \cdot & \cdot & 0 & 0 & 0 & 0 \\ \rho & \rho & 1 & \rho & \cdot & \cdot & \cdot & \cdot & 0 & 0 & 0 & 0 \\ \rho & \rho & \rho & 1 & \cdot & \cdot & \cdot & \cdot & 0 & 0 & 0 & 0 \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ 0 & 0 & 0 & 0 & \cdot & \cdot & \cdot & \cdot & 1 & \rho & \rho & \rho \\ 0 & 0 & 0 & 0 & \cdot & \cdot & \cdot & \cdot & \rho & 1 & \rho & \rho \\ 0 & 0 & 0 & 0 & \cdot & \cdot & \cdot & \cdot & \rho & \rho & 1 & \rho \\ 0 & 0 & 0 & 0 & \cdot & \cdot & \cdot & \cdot & \rho & \rho & \rho & 1 \end{bmatrix}$$

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

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

31

Here we have a total of 3 regression parameters and variance components  $(\beta_0, \sigma_0, \sigma_\epsilon)$ , 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.

32

## Inference for Regression Parameters

The score equation for  $\boldsymbol{\beta}$  is

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

and yields the MLE for  $\boldsymbol{\beta}$  as

$$\hat{\boldsymbol{\beta}} = \left( \sum_{i=1}^m \mathbf{x}_i^T \mathbf{V}_i^{-1} \mathbf{x}_i \right)^{-1} \left( \sum_{i=1}^m \mathbf{x}_i^T \mathbf{V}_i^{-1} \mathbf{y}_i \right), \quad (13)$$

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

The variance of  $\hat{\boldsymbol{\beta}}$  may be obtained either directly from (13), or from the second derivative of the log-likelihood. Since

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

the observed and expected information matrices coincide with

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

33

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

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

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

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

where  $\hat{\boldsymbol{\alpha}}$  is a consistent estimator of  $\boldsymbol{\alpha}$ . 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  $\boldsymbol{\alpha}$ , see Verbeke and Molenberghs (2000, Chapter 6), but if the sampling size is not sufficiently large for reliable estimation of  $\boldsymbol{\alpha}$ , we recommend following a Bayesian approach to inference.

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

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

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

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

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  $\boldsymbol{\beta}$  and  $\boldsymbol{\alpha}$ . Secondly, asymptotically we have independence between  $\hat{\boldsymbol{\beta}}$  and  $\hat{\boldsymbol{\alpha}}$ , so any consistent estimator of  $\boldsymbol{\alpha}$  will give an asymptotically efficient estimator for  $\boldsymbol{\beta}$ .

Likelihood ratio tests are available for regression parameters.

35

### Inference for Variance Components by MLE

The MLE of  $\boldsymbol{\alpha}$  follows from maximization of (12), 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  $\boldsymbol{\beta}$ .

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

An alternative and often preferable method is provided by restricted maximum likelihood.



## 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

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

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  $\sigma_0^2$  is a non-negative scalar. In this case the asymptotic null distribution is a 50:50 mixture of  $\chi_0^2$  and  $\chi_1^2$  distributions, where the former is the distribution that gives probability mass 1 to the value 0.

If the usual  $\chi_1^2$  distribution is used then the null would be accepted too often, leading to a variance component structure that is too simple.

Estimating  $\sigma_0^2$  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.

37

## Inference for Variance Components by REML

Restricted 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.

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

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

39

### 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

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

## 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  $\boldsymbol{\beta}$ .

41

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

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

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

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

Then

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

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

Further  $\mathbf{B}^\top \mathbf{X} = \mathbf{0}$ , so that

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

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  $\mathbf{Y} \rightarrow (\mathbf{U}, \hat{\boldsymbol{\beta}}_G) = (\mathbf{B}^\top \mathbf{Y}, \mathbf{G}^\top \mathbf{Y})$ , where

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

is the generalized least squares estimator.

42

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.  $\begin{vmatrix} \mathbf{T} & \mathbf{U} \\ \mathbf{V} & \mathbf{W} \end{vmatrix} = |\mathbf{T}| |\mathbf{W} - \mathbf{V} \mathbf{T}^{-1} \mathbf{U}|$ .

Then

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

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

43

We have

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

and

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

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

Hence

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

Let  $\mathbf{S}_1, \mathbf{S}_2$ , be minimal sufficient statistics for which

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

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

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  $\mathbf{s}_1 = \mathbf{u}$ ,  $\mathbf{s}_2 = \widehat{\boldsymbol{\beta}}_G$ ,  $\boldsymbol{\lambda} = \boldsymbol{\alpha}$ ,  $\boldsymbol{\phi} = \boldsymbol{\beta}$ , and  $p(\mathbf{U} \mid \boldsymbol{\alpha})$  is a marginal likelihood.

Hence

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

We have

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

and

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

This leads to

$$\begin{aligned} p(\mathbf{U} \mid \boldsymbol{\alpha}) &= (2\pi)^{-(N-k-1)/2} \frac{|\mathbf{x}^T \mathbf{x}|^{1/2} |\mathbf{V}|^{-1/2}}{|\mathbf{x}^T \mathbf{V}^{-1} \mathbf{x}|^{1/2}} \\ &\times \exp \left\{ -\frac{1}{2} (\mathbf{y} - \mathbf{x}\widehat{\boldsymbol{\beta}}_G)^T \mathbf{V}^{-1} (\mathbf{y} - \mathbf{x}\widehat{\boldsymbol{\beta}}_G) \right\} \end{aligned} \quad (14)$$

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

45

- 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,

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

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

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

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.

46

- 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.

47

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  $\boldsymbol{\epsilon}_i$ .

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

$$\hat{\sigma}^2 = \frac{\sum_{i=1}^m \boldsymbol{\epsilon}_i^T \boldsymbol{\epsilon}_i}{\sum_{i=1}^m n_i} = \frac{t_1}{N}$$

$$\hat{\mathbf{D}} = \frac{1}{m} \sum_{i=1}^m \mathbf{b}_i \mathbf{b}_i^T = \frac{\mathbf{t}_2}{m},$$

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:

$$\hat{t}_1 = \mathbb{E} \left[ \sum_{i=1}^m \boldsymbol{\epsilon}_i^T \boldsymbol{\epsilon}_i | \mathbf{y}_i, \hat{\beta}(\hat{\alpha}), \hat{\alpha} \right]$$

$$\hat{\mathbf{t}}_2 = \mathbb{E} \left[ \sum_{i=1}^m \mathbf{b}_i^T \mathbf{b}_i | \mathbf{y}_i, \hat{\beta}(\hat{\alpha}), \hat{\alpha} \right].$$

Closed form fixed and random effect estimates are available once we know  $\boldsymbol{\alpha}$ .

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

Let  $\boldsymbol{\theta}$  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  $\boldsymbol{\theta}^{(t)}$ , the estimate at iteration  $t$  gives:

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

differentiating and setting equal to zero:

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

gives the next estimate

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

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.

49

## Dental Example

The simplest possible mixed effects model is given by

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

where  $\epsilon_{ij}$  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  $\beta_0$ .

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.

The marginal distribution is normal with mean

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

where

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

is  $4m \times 1$  vector and

$$\boldsymbol{\mu}_i = (\beta_0 + \beta_1 t_1, \beta_0 + \beta_1 t_2, \beta_0 + \beta_1 t_3, \beta_0 + \beta_1 t_4)^\top.$$

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

51

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

52



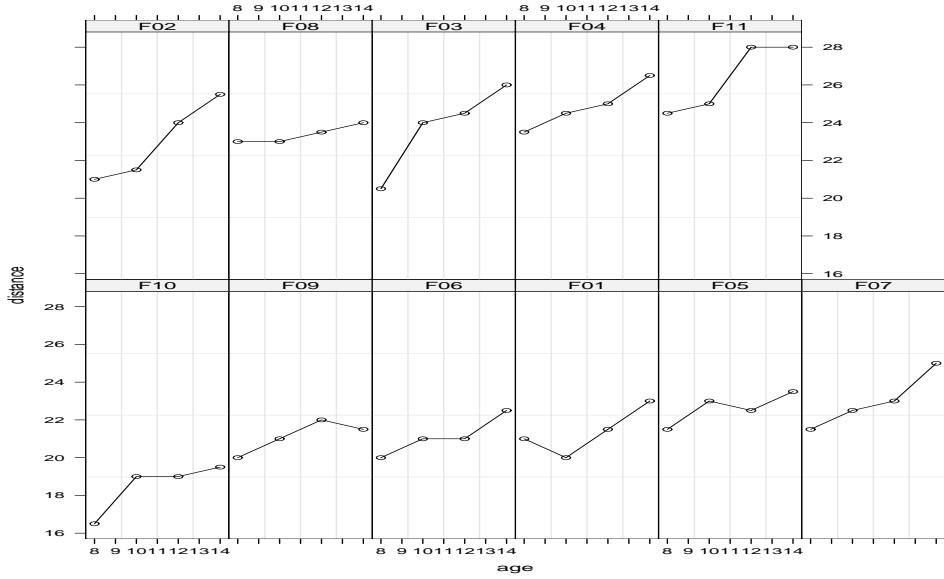


Figure 6: Length versus age (in years) for 11 girls.

53

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
```

54

Notice the standard error for  $\beta_1$  is smaller for the REML analysis – slopes are being estimated from within-girl comparisons (for the intercept correct standard errors are larger due to dependence).

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 (10). 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 (10).

55

Now for comparison carry out 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).

56

## 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:

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

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

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

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

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

57

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  $p(\boldsymbol{\theta} | \mathbf{y}) = 0$  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*.

## Sequential Updating

Suppose first that  $\mathbf{y}_1$  and  $\mathbf{y}_2$  represent the current totality of data. Then the posterior is given by

$$p(\boldsymbol{\theta} \mid \mathbf{y}_1, \mathbf{y}_2) = \frac{p(\mathbf{y}_1, \mathbf{y}_2 \mid \boldsymbol{\theta})\pi(\boldsymbol{\theta})}{p(\mathbf{y}_1, \mathbf{y}_2)}. \quad (15)$$

Now suppose that we are at a previous time point at which only  $\mathbf{y}_1$  are available, the posterior in this case is

$$p(\boldsymbol{\theta} \mid \mathbf{y}_1) = \frac{p(\mathbf{y}_1 \mid \boldsymbol{\theta})\pi(\boldsymbol{\theta})}{p(\mathbf{y}_1)}.$$

When  $\mathbf{y}_2$  becomes available, the “prior” for these data corresponds to  $p(\boldsymbol{\theta} \mid \mathbf{y}_1)$  since it represents the current beliefs concerning  $\boldsymbol{\theta}$ . We then update via

$$p(\boldsymbol{\theta} \mid \mathbf{y}_1, \mathbf{y}_2) = \frac{p(\mathbf{y}_2 \mid \mathbf{y}_1, \boldsymbol{\theta})\pi(\boldsymbol{\theta} \mid \mathbf{y}_1)}{p(\mathbf{y}_2 \mid \mathbf{y}_1)}. \quad (16)$$

Identical inference in each case; hence consistent inference is reached regardless of whether we produce the posterior in one stage or two, corresponding to whether all of the data are analyzed simultaneously.

59

## 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  $\theta_i$  is given by

$$p(\theta_i \mid \mathbf{y}) = \int_{\boldsymbol{\theta}_{-i}} p(\boldsymbol{\theta} \mid \mathbf{y}) \, d\boldsymbol{\theta}_{-i}, \quad (17)$$

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

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

$$E[\theta_i \mid \mathbf{y}] = \int_{\theta_i} \theta_i p(\theta_i \mid \mathbf{y}) \, d\theta_i. \quad (18)$$

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

$$\int_{-\infty}^{\theta_i(q)} p(\theta_i \mid \mathbf{y}) \, d\theta_i. \quad (19)$$

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

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

$$p(\mathbf{z} | \mathbf{y}) = \int_{\boldsymbol{\theta}} p(\mathbf{z} | \boldsymbol{\theta})p(\boldsymbol{\theta} | \mathbf{y}) d\boldsymbol{\theta}. \quad (20)$$

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, (17)–(20), given that for most models, these are analytically intractable

61

**Simple Example:** Suppose we have

$$Y_i | \theta \sim_{i.i.d.} N(\theta, \sigma^2), \quad i = 1, \dots, n,$$

with  $\sigma^2$  assumed known and  $\theta$  unknown.

*Estimation*

Recall that the MLE

$$\bar{Y} \sim N\left(\theta, \frac{\sigma^2}{n}\right).$$

Suppose the prior distribution for  $\theta$  can be described by a normal distribution with mean  $m$  and variance  $v$  ( $m$  and  $v$  are known). Then the posterior distribution  $p(\theta | \mathbf{y})$  is given by

$$N\left(\bar{y} \times w + m \times (1 - w), \frac{\sigma^2}{n} \times w\right),$$

where  $w = \frac{v}{v + \sigma^2/n}$ .

Think about cases:  $n = 0$  (recover the prior),  $v = 0$  (posterior=prior),  $v^{-1} = 0$  (improper prior, frequentist and Bayesian estimates coincide),  $n \rightarrow \infty$  ( $w \rightarrow 1$  unless  $v = 0$ ).

62

One useful way of specifying the prior is as

$$\theta \sim N\left(m, \frac{\sigma^2}{k}\right),$$

in which case  $k$  may be regarded as a *prior sample size*. It is ‘as if’ we carried out an experiment with  $k$  observations and we observed a mean of  $m$ . This gives  $w = n/(n + k)$ .

63

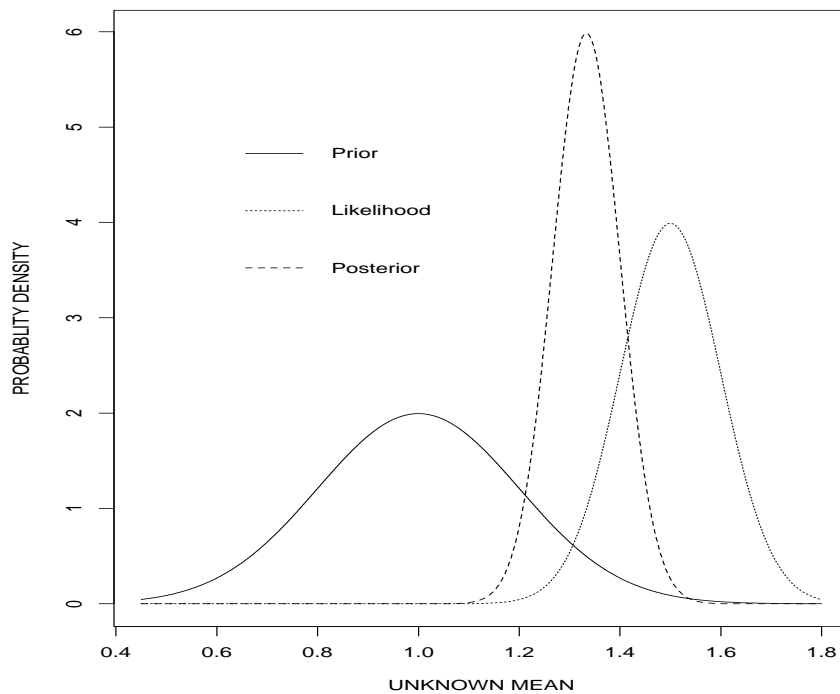


Figure 7: Normal likelihood ( $\bar{y}=1.5, n=10, \sigma=1$ ), normal prior ( $m=1, k=5$ ) and the resultant normal posterior.

64

### Prediction

Suppose we wish to obtain the predictive density for a new random variable  $Z \sim N(\theta, \sigma^2)$ .

Then

$$p(z|\mathbf{y}) = \int p(z|\theta) \times p(\theta|\mathbf{y})d\theta.$$

It may be shown that

$$z|\mathbf{y} \sim N \{E[\theta|\mathbf{y}], \sigma^2 + \text{var}(\theta|\mathbf{y})\},$$

so that the mean of the predictive distribution is the posterior mean and the variance is given by the sum of the ‘measurement error’ and the uncertainty in the posterior mean.

65

### Prior Choice

We distinguish between two prior specification situations. In the first, which we label as a *baseline prior* an analysis is required in which the prior distribution has minimal impact, so that the information in the likelihood dominates the posterior.

The second situation, which we label as a *substantive prior* is one in which it is desired to incorporate more substantial prior information into the analysis.

#### Baseline Priors

On first consideration it would seem that the specification of a baseline prior is straightforward, one simply takes the choice

$$\pi(\boldsymbol{\theta}) \propto 1 \tag{21}$$

so that the posterior distribution depends solely on the data through the likelihood  $p(\mathbf{y} | \boldsymbol{\theta})$ .

There are two difficulties with this.

66

The first difficulty is that the prior (21) is improper (it does not integrate to a positive constant  $< \infty$ ) unless the range of each element of  $\boldsymbol{\theta}$  is finite.

In some instances this is not a problem since the posterior corresponding to the prior is proper. Philosophically a posterior arising from an improper prior may be justified as a limiting case of proper priors. More practically we may instead assume that the prior is integrable over its support but is “locally uniform”, so that the likelihood dominates.

For nonlinear models in particular, care must be taken to ensure that the posterior corresponding to a particular prior choice is proper. Some general guidelines are available, for example, improper priors for the regression parameters in a generalized linear model will usually lead to a proper posterior although not for some pathological cases.

For example suppose  $Y | p \sim \text{Binomial}(n, p)$ , and a uniform prior is used on the logit of  $p$ ,  $\log\{p/(1-p)\}$  which implies the prior on  $p$  is  $\pi(p) = [p(1-p)]^{-1}$ . Then an improper posterior results if  $y = 0$  (or  $y = n$ ) since the non-integrable spike at  $p = 0$  (or  $p = 1$ ) remains in the posterior. For  $n = 1$  one of these events will always occur and so an improper posterior always results.

67

To illustrate the non propriety in another non-linear situation consider the model

$$Y_i | \theta \sim_{ind} N\{\exp(-\theta x_i), \sigma^2\}, \quad (22)$$

$i = 1, \dots, n$ , with  $\theta > 0$  and  $\sigma^2$  assumed known. With an improper uniform prior on  $\theta$  we have the posterior

$$p(\theta | \mathbf{y}) \propto \exp \left\{ -\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - e^{-\theta x_i})^2 \right\}.$$

As  $\theta \rightarrow \infty$ ,

$$p(\theta | \mathbf{y}) \rightarrow \exp \left\{ -\frac{1}{2\sigma^2} \sum_{i=1}^n y_i^2 \right\},$$

a constant, so that the posterior is improper.

68



The second is that if we reparameterize the model in terms of  $\boldsymbol{\phi} = \mathbf{g}(\boldsymbol{\theta})$  where  $\mathbf{g}(\cdot)$  is a one-one mapping, then the prior for  $\boldsymbol{\phi}$  corresponding to (21) is given by

$$\pi(\boldsymbol{\phi}) = \left| \frac{d\boldsymbol{\theta}}{d\boldsymbol{\phi}} \right|,$$

which, unless  $\mathbf{g}$  is linear, is not constant.

As an example, consider a variance  $\sigma^2$ , the prior  $\pi(\sigma^2) \propto 1$  corresponds to a prior for the standard deviation of  $\pi(\sigma) \propto \sigma$ ; the problem is that we cannot be “flat” on different scales.

This indicates that a desirable property in constructing baseline priors is there invariance to parameterization, so that we obtain the same prior regardless of the starting parameterization. In the example just considered suppose the data are normally distributed with variance  $\sigma^2$ . The improper prior

$$\pi(\sigma) \propto \frac{1}{\sigma}$$

has a number of justifications including invariance to parameterization.

69

### 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

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

70

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

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

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

71

Hence the posterior

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

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

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

72

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

73

Hence

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

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

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

74

**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

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

which differs from the “usual” likelihood by the term

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

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  $\hat{\boldsymbol{\alpha}}$ .

75

**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  $\mathbf{b}$ .

76

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

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

77

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:  $E_u[U] = E_{u,v}[U] = E_v\{E_{u|v}[U|V]\}$ .)

78

## 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

$$\begin{aligned} \text{cov}_{\tilde{\mathbf{b}}, \mathbf{b}}(\tilde{\mathbf{b}}_i, \mathbf{b}_i) &= E_{\mathbf{y}}[\text{cov}(\tilde{\mathbf{b}}_i, \mathbf{b}_i \mid \mathbf{y})] + \text{cov}_{\mathbf{y}}(E[\tilde{\mathbf{b}}_i \mid \mathbf{y}], E[\mathbf{b}_i \mid \mathbf{y}]) \\ &= E_{\mathbf{y}}[\text{cov}(\tilde{\mathbf{b}}_i, \mathbf{b}_i \mid \mathbf{y})] + \text{cov}_{\mathbf{y}}(\tilde{\mathbf{b}}_i, \tilde{\mathbf{b}}_i) \\ &= \text{var}(\tilde{\mathbf{b}}_i) \end{aligned} \quad (23)$$

The first term in (23) is the covariance between a constant  $E[\tilde{\mathbf{b}} \mid \mathbf{y}]$  (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  $E[\tilde{\mathbf{b}}_i \mid \mathbf{y}] = E[E[\mathbf{b}_i \mid \mathbf{y}] \mid \mathbf{y}] = \tilde{\mathbf{b}}_i$ .

Hence

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

79

## Application to the LMEM

The predictor,  $\tilde{\mathbf{b}} = E[\mathbf{b} \mid \mathbf{y}]$ , is a random variable, since it a function of  $\mathbf{y}$ , and so we need to know something about  $p(\mathbf{b} \mid \mathbf{y})$  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  $\mathbf{V} = \mathbf{A}\mathbf{U}$  where  $\mathbf{A}$  is an  $m \times n$  matrix. Then  $\text{cov}(\mathbf{U}, \mathbf{A}\mathbf{U}) = \mathbf{W}\mathbf{A}^T$  where  $\mathbf{W} = \text{cov}(\mathbf{U})$ , and  $\text{cov}(\mathbf{A}\mathbf{U}, \mathbf{U}) = \mathbf{A}\mathbf{W}$ .

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  $\mathbf{b}_i \sim N(\mathbf{0}, \mathbf{D})$ ,  $\boldsymbol{\epsilon} \sim N(\mathbf{0}, \sigma_{\boldsymbol{\epsilon}}^2 \mathbf{I})$  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  $\text{cov}(\mathbf{y}_i, \mathbf{b}_i) = \mathbf{z}_i \mathbf{D}$ .

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

$$\tilde{\mathbf{b}}_i = E[\mathbf{b}_i | \mathbf{y}_i] = \mathbf{D} \mathbf{z}_i^T \mathbf{V}_i^{-1} (\mathbf{y}_i - \mathbf{x}_i \boldsymbol{\beta}) \quad (24)$$

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 data towards zero, as we see in examples later.

The form (24) is not of practical use since it depends on the unknown  $\boldsymbol{\beta}$  and  $\boldsymbol{\alpha}$ ; instead we use

$$\tilde{\mathbf{b}}_i = E[\mathbf{b}_i | \mathbf{y}_i] = \hat{\mathbf{D}} \mathbf{z}_i^T \hat{\mathbf{V}}_i^{-1} (\mathbf{y}_i - \mathbf{x}_i \hat{\boldsymbol{\beta}}). \quad (25)$$

Substitution of  $\hat{\boldsymbol{\beta}}$  is not such a problem (since it is an unbiased estimator, and appears in (24) in a linear fashion), but  $\hat{\boldsymbol{\alpha}}$  is more problematic.

81

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  $\hat{\boldsymbol{\alpha}}$ .

82

We now examine fitted values:

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

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

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

We can also write

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

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

83

### Example: One-way ANOVA

For the simple balanced ANOVA model previously considered

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

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

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

$$\widehat{\beta}_0 + \widetilde{b}_i \rightarrow \bar{y}_i \rightarrow \mathbf{E}[Y_i].$$

84



The form of (24) 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

$$\tilde{\mathbf{b}}_i = \mathbf{D}\mathbf{z}_i^T \mathbf{V}_i^{-1}(\mathbf{y}_i - \mathbf{x}_i\boldsymbol{\beta}).$$

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.

85

### 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
> coef(remlelm)
      (Intercept) I(age - 11)
F10    18.64240    0.4795455
F09    21.17728    0.4795455
F06    21.17728    0.4795455
F01    21.41869    0.4795455
F05    22.62578    0.4795455
F07    22.98791    0.4795455
F02    22.98791    0.4795455
F08    23.35003    0.4795455
F03    23.71216    0.4795455
F04    24.79853    0.4795455
F11    26.24704    0.4795455
```

86

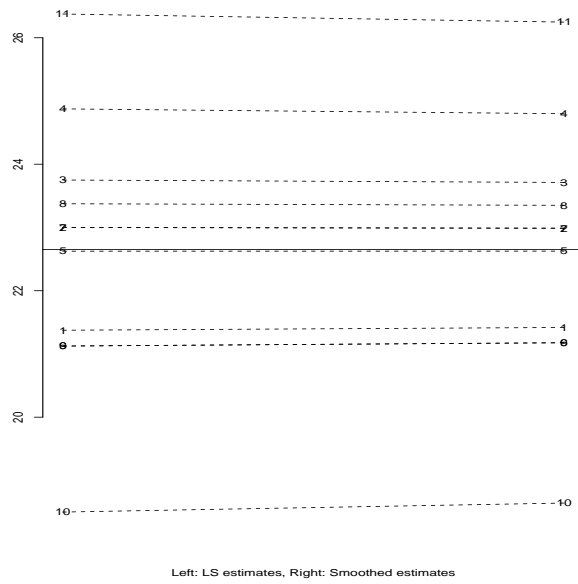


Figure 8: Least squares estimates and smoothed estimates,  $\hat{\beta}_0 + \tilde{b}_i$ ; not much shrinkage here since  $\hat{\sigma}_0$  is large relative to  $\hat{\sigma}_\epsilon$ .

87

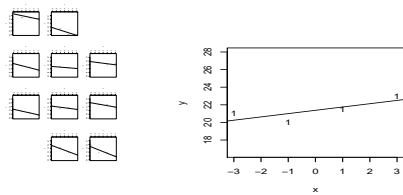


Figure 9: Individual fits: solid line is LS, broken line from LMEM.

88

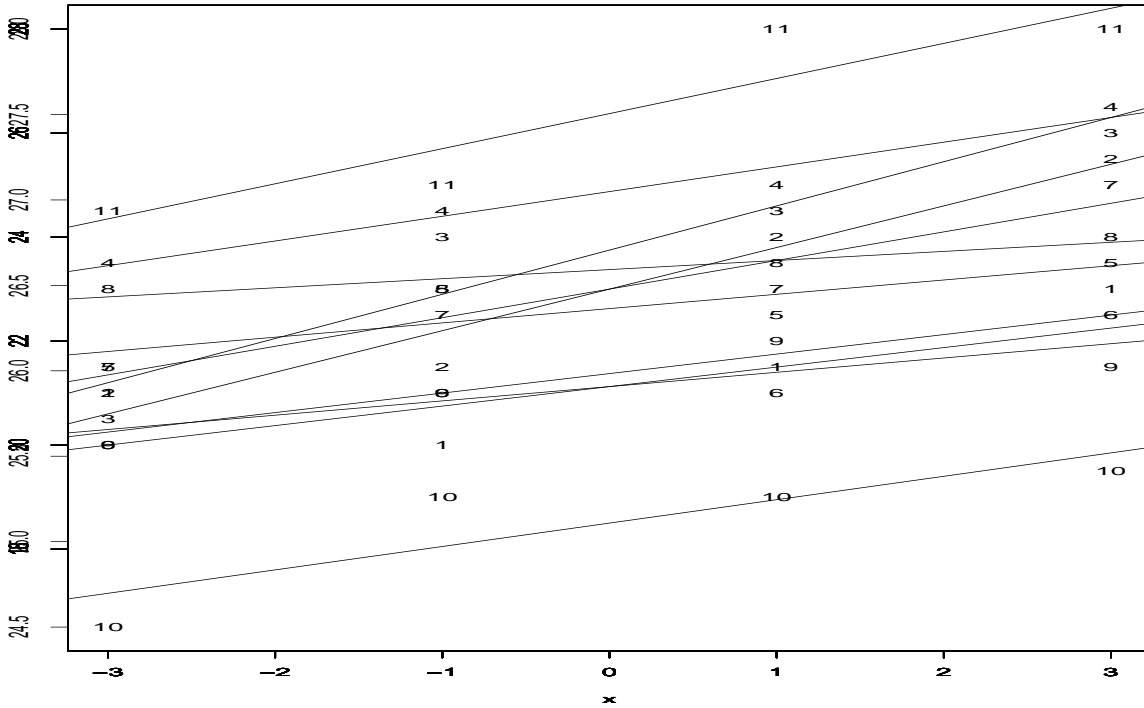


Figure 10: Individual fits: Solid lines are LS, broken line from LMEM.

Dental Example: Boys and Girls Joint Analyses

Table 2 describes a variety of LMEMs 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) )

```

91

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

Figure 12 gives a scatter plot of the LS estimates of intercepts and slopes, for boys and girls.

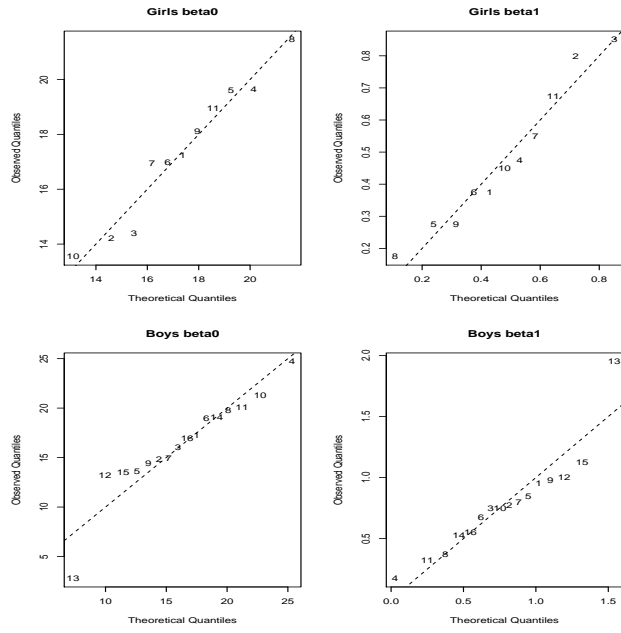


Figure 11: QQ plot of the LS estimates.

93

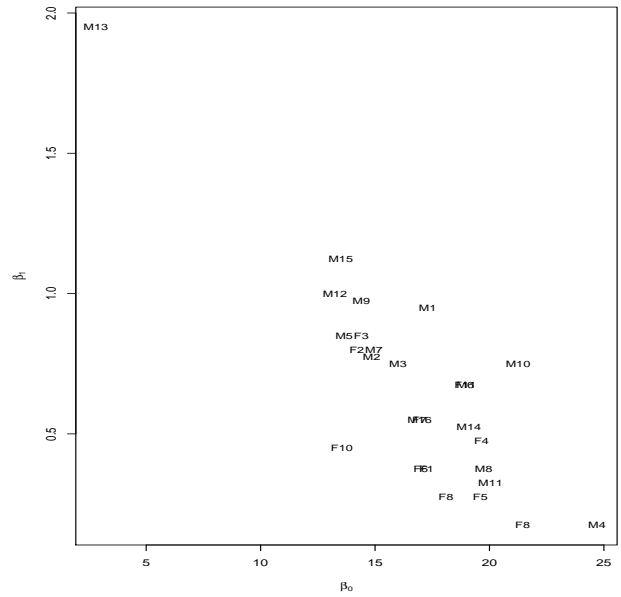


Figure 12: Plot of the LS estimates for boys and girls.

94

## Bayesian Inference for the LMEM

Consider the model

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

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

The form of the posterior follows from exploiting conditional independencies:

$$\begin{aligned} p(\boldsymbol{\beta}, \boldsymbol{\alpha}, \mathbf{b} \mid \mathbf{y}) &\propto p(\mathbf{y} \mid \boldsymbol{\beta}, \boldsymbol{\alpha}, \mathbf{b})\pi(\boldsymbol{\beta}, \boldsymbol{\alpha}, \mathbf{b}) = \prod_{i=1}^m p(\mathbf{y}_i \mid \boldsymbol{\beta}, \boldsymbol{\alpha}, \mathbf{b}_i)\pi(\mathbf{b} \mid \boldsymbol{\alpha})\pi(\boldsymbol{\beta})\pi(\boldsymbol{\alpha}) \\ &= \prod_{i=1}^m \{p(\mathbf{y}_i \mid \boldsymbol{\beta}, \boldsymbol{\alpha}, \mathbf{b}_i)\pi(\mathbf{b}_i \mid \boldsymbol{\alpha})\} \pi(\boldsymbol{\beta})\pi(\boldsymbol{\alpha}) \end{aligned} \quad (26)$$

Alternatively, we can derive the posterior for  $\boldsymbol{\beta}, \boldsymbol{\alpha}$  directly:

$$\begin{aligned} p(\boldsymbol{\beta}, \boldsymbol{\alpha} \mid \mathbf{y}) &\propto p(\mathbf{y} \mid \boldsymbol{\beta}, \boldsymbol{\alpha})\pi(\boldsymbol{\beta}, \boldsymbol{\alpha}) = \prod_{i=1}^m p(\mathbf{y}_i \mid \boldsymbol{\beta}, \boldsymbol{\alpha})\pi(\boldsymbol{\beta}, \boldsymbol{\alpha}) \\ &= \prod_{i=1}^m \int p(\mathbf{y}_i, \mathbf{b}_i \mid \boldsymbol{\beta}, \boldsymbol{\alpha}) d\mathbf{b}_i \pi(\boldsymbol{\beta}, \boldsymbol{\alpha}) \end{aligned}$$

where the integrand is given by the term in curly brackets in (26).

The prior on  $\mathbf{b}_i$  is justified by the context, formally via *exchangeability*.

95

## Exchangeability

**Definition:** A finite set  $Y_1, \dots, Y_n$  of random variables is said to be *exchangeable* if every permutation  $(Y_1, \dots, Y_n)$  has the same joint distribution as every other permutation. An infinite collection is exchangeable if every finite subcollection is exchangeable.

Every collection of independent and identically distributed random variables is exchangeable.

**Theorem:** *De Finetti's representation Theorem for 0/1 random variables.*

If  $Y_1, Y_2, \dots$  is an infinitely exchangeable sequence of 0/1 random variables, there exists a distribution  $\pi(\cdot)$  such that the joint mass function  $\Pr(y_1, \dots, y_n)$  has the form

$$\Pr(y_1, \dots, y_n) = \int_0^1 \prod_{i=1}^n \theta^{y_i} (1 - \theta)^{1-y_i} \pi(\theta) d\theta,$$

where

$$\int_0^\theta \pi(u) du = \lim_{n \rightarrow \infty} \Pr\left(\frac{Z_n}{n} \leq \theta\right),$$

with  $Z_n = Y_1 + \dots + Y_n$ , and  $\theta = \lim_{n \rightarrow \infty} Z_n/n$ .

96

**Proof:** See Bernardo and Smith (1994) for more details.

Let  $z_n = y_1 + \dots + y_n$  be the number of 1's (which we label "successes") in the first  $n$  observations. Then, due to exchangeability,

$$\Pr(y_1 + \dots + y_n = z_n) = \binom{n}{z_n} \Pr(Y_{\pi(1)}, \dots, Y_{\pi(n)}),$$

for all permutations  $\pi$  of  $\{1, \dots, n\}$  such that  $y_{\pi(1)} + \dots + y_{\pi(n)} = z_n$ . Then we can embed the event  $y_1 + \dots + y_n = z_n$  within a sequence,  $y_1, \dots, y_N$ ,  $N \geq n$ , and

$$\begin{aligned} \Pr\left(\sum_{i=1}^n y_i = z_n\right) &= \sum_{z_N=z_n}^{N-(n-z_n)} \Pr(y_1 + \dots + y_n = z_n, y_1 + \dots + y_N = z_N) \\ &= \sum_{z_N=z_n}^{N-(n-z_n)} \Pr(y_1 + \dots + y_n = z_n \mid y_1 + \dots + y_N = z_N) \\ &\quad \times \Pr(y_1 + \dots + y_N = z_N). \end{aligned}$$

To obtain the conditional probability we observe that it is as if we have a population of  $N$  people of which  $z_N$  are successes, and  $N - z_N$  failures, from which we draw  $n$  people, the probability of  $z_n$  successes is then hypergeometric.

97

Hence

$$\Pr(y_1 + \dots + y_n = z_n) = \sum_{z_N=z_n}^{N-(n-z_n)} \frac{\binom{z_N}{z_n} \binom{N-z_N}{n-z_n}}{\binom{N}{n}} \Pr(z_N)$$

Here  $\Pr(z_N)$  is the "prior" belief in the number of successes out of  $N$ .

Let  $N \rightarrow \infty$  and by the strong law of law numbers  $\theta = \lim_{N \rightarrow \infty} z_N/N$ .

The hypergeometric tends to a binomial with parameters  $n$  and  $\theta$ , and the prior  $\Pr(z_N)$  is translated into a prior for  $\theta$ ,  $\pi(\theta)$ . Hence we have

$$\Pr(y_1 + \dots + y_n = z_n) \rightarrow \binom{n}{z_n} \int \theta^{z_n} (1-\theta)^{n-z_n} \pi(\theta) d\theta,$$

as  $N \rightarrow \infty$ .

98

## Implications

The interpretation of this theorem is of great significance:

- We may view the  $Y_i$  to be independent, Bernoulli random variables, conditional on a random variable  $\theta$ .
- $\theta$  is itself assigned a probability distribution  $\pi(\cdot)$ .
- $\pi$  may be interpreted as ‘beliefs about the limiting relative frequency of 1’s’.

In conventional language, we have the *likelihood function*

$$p(Y_1, \dots, Y_n | \theta) = \prod_{i=1}^n p(Y_i | \theta) = \prod_{i=1}^n \theta^{Y_i} (1 - \theta)^{1 - Y_i},$$

where the *parameter*  $\theta$  is assigned a *prior distribution*  $\pi(\theta)$ .

99

**Corollary:** If  $Y_1, Y_2, \dots$  is an infinitely exchangeable sequence of 0/1 random variables, then we have the conditional probability function

$$p(y_{m+1}, \dots, y_n | y_1, \dots, y_m) = \int_0^1 \prod_{i=m+1}^n \theta^{Y_i} (1 - \theta)^{1 - Y_i} \pi(\theta | y_1, \dots, y_m) d\theta,$$

for  $1 \leq m < n$  where

$$\pi(\theta | y_1, \dots, y_m) = \frac{\prod_{i=1}^m \theta^{y_i} (1 - \theta)^{1 - y_i} \pi(\theta)}{\int_0^1 \prod_{i=1}^m \theta^{y_i} (1 - \theta)^{1 - y_i} \pi(\theta) d\theta}$$

and

$$\int_0^\theta \pi(u) du = \lim_{n \rightarrow \infty} \Pr\left(\frac{z_n}{n} \leq \theta\right).$$

## Proof

Write

$$\Pr(y_{m+1}, \dots, y_n | y_1, \dots, y_m) = \frac{\Pr(y_1, \dots, y_n)}{\Pr(y_1, \dots, y_m)},$$

and then use the previous result on numerator and denominator.

**Interpretation:** the *prior distribution*  $\pi(\theta)$  for  $\theta$  has been revised, via *Bayes’ Theorem*, into the *posterior distribution*  $\pi(\theta | y_1, \dots, y_m)$ .



## Further results

### General Representation Theorem:

If  $Y_1, Y_2, \dots$  is an infinitely exchangeable sequence of random variables with probability measure  $P$ , there exists a distribution function  $Q$  such that the joint mass function  $p(Y_1, \dots, Y_n)$  has the form

$$p(Y_1, \dots, Y_n) = \int \prod_{i=1}^n p(Y_i | \boldsymbol{\theta}) \pi(\boldsymbol{\theta}) d\boldsymbol{\theta},$$

with  $p(\cdot | \boldsymbol{\theta})$  denoting the density function corresponding to the ‘unknown parameter’  $\boldsymbol{\theta}$ .

Further assumptions on  $Y_1, Y_2, \dots$  are required to identify  $p(\cdot | \boldsymbol{\theta})$ .

101

## Relevance of Exchangeability

If we believe *a priori* that  $\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_m$  are exchangeable (and are considered within a hypothetical infinite sequence of such random variables), then it can be shown using representation theorems that the prior can be written in the form

$$p(\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_m) = \int \prod_{i=1}^m p(\boldsymbol{\theta}_i | \boldsymbol{\phi}) \pi(\boldsymbol{\phi}) d\boldsymbol{\phi},$$

that is, they are conditionally independent, given *hyperparameters*  $\boldsymbol{\phi}$ , with the hyperparameters having a *hyperprior* distribution.

Hence we have a two-stage (hierarchical) prior:

*Stage A:*  $\boldsymbol{\theta}_i | \boldsymbol{\phi} \sim_{iid} p(\cdot | \boldsymbol{\phi})$ ,  $i = 1, \dots, m$ .

*Stage B:*  $\boldsymbol{\phi} \sim_{iid} \pi(\cdot)$ .

Parametric choices for  $p(\cdot | \boldsymbol{\phi})$  and  $\pi(\cdot)$  are usually made for computational convenience.

Contrast with the sampling theory approach in which the random effects are assumed to be a random sample from a hypothetical infinite population.

102

## Bayesian Computation

We have seen that to summarize posterior distributions integration is required and, in all but the simplest (conjugate) models, these integrals are not analytically tractable.

Integration is also required to integrate out the random effects in nonlinear mixed effects models, to obtain the likelihood, and later we will review a number of analytical and numerical approaches, for now we concentrate on Markov chain Monte Carlo (MCMC).

The first key idea is the duality between densities and samples from that density: given a density we can always generate samples, and given samples we can reconstruct the density.

Simulation-based techniques have revolutionized Bayesian statistics, by allowing the fitting of very complex models.

103

## Example

Suppose we have  $Y_j | p_j \sim \text{Binomial}(n_j, p_j)$ ,  $j = 1, 2$ , with independent priors  $p_j \sim U(0, 1)$ . The posteriors are available analytically as  $p_j | y_j \sim \text{Beta}(y_j + 1, n_j - y_j + 1)$ , but suppose we are interested in inference for the odds ratio  $\phi = \frac{p_1}{1-p_1} / \frac{p_2}{1-p_2}$  and for the relative risk  $\theta = \frac{p_1}{p_2}$ .

The following is R code to simulate from  $\phi | y_1, y_2$  when  $n_1 = 35, n_2 = 45, y_1 = 30, y_2 = 10$ :

```
> n1 <- 35; n2 <- 45; y1 <- 30; y2 <- 10
> nsamp <- 1000
> p1 <- rbeta(nsamp, y1+1, n1-y1+1); p2 <- rbeta(nsamp, y2+1, n2-y2+1)
> odds <- (p1/(1-p1))/(p2/(1-p2)); rr <- p1/p2
> par(mfrow=c(2,2))
> hist(p1, xlim=c(0,1))
> hist(p2, xlim=c(0,1))
> hist(odds)
> hist(rr)
> sum(odds[odds>10])/sum(odds) # Posterior prob that odds ratio is > than 10
[1] 0.945683
```

104

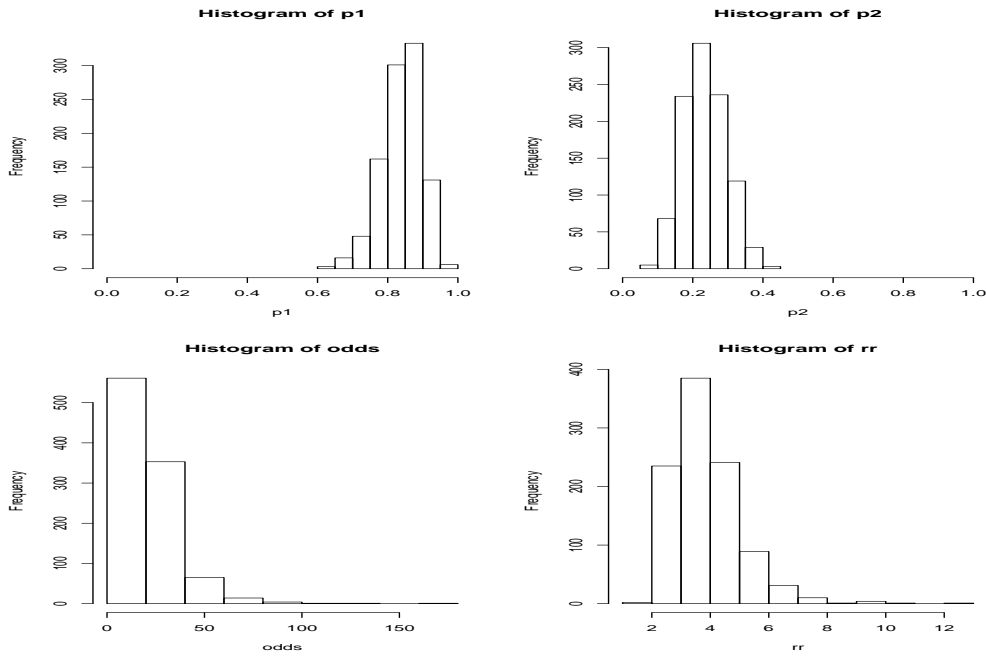


Figure 13: Posterior distributions for  $p_1$ ,  $p_2$ , the odds ratio  $\frac{p_1}{1-p_1} / \frac{p_2}{1-p_2}$  and for the relative risk  $\theta = \frac{p_1}{p_2}$ .

105

### The Composition Method

A useful technical for simulating from joint posterior distributions is the following.

Write the joint posterior distribution for  $\theta_1, \theta_2$  as

$$p(\theta_1, \theta_2 | \mathbf{y}) = p(\theta_1 | \mathbf{y})p(\theta_2 | \theta_1, \mathbf{y})$$

Then a simulating algorithm to produce independent samples from  $p(\theta_1, \theta_2 | \mathbf{y})$  is, for  $s = 1, \dots, S$ :

1. Simulate  $\theta_1^{(s)} \sim_{ind} p(\theta_1 | \mathbf{y})$ .
2. Simulate  $\theta_2^{(s)} \sim_{ind} p(\theta_2 | \theta_1^{(s)}, \mathbf{y})$ .

## Markov chain Monte Carlo

MCMC is a very general technique that has revolutionized practical Bayesian statistics.

In the usual derivation of Markov chains over a discrete sample space we are given a transition matrix and the aim is to find the stationary distribution (if it exists). Probabilities of movement depend on the current state only, hence the name.

In the context of sampling from a distribution  $\pi(\cdot)$ , the aim is to construct a Markov chain whose stationary distribution is  $\pi$ .

Samples  $\boldsymbol{\theta}^{(s)}$ ,  $s = 1, \dots, S$ , produced by a Markov chain “look” more and more like *dependent* samples from  $\pi$  as  $S \rightarrow \infty$ . The dependency does not cause a problem in terms of estimation since

$$\frac{1}{S} \sum_{s=1}^S f(\boldsymbol{\theta}^{(s)}) \rightarrow \text{E}\{f(\boldsymbol{\theta})\},$$

as  $S \rightarrow \infty$  (provided the expectation exists).

The only difficulty with the dependency is establishing an appropriate Monte Carlo error on the resultant estimator. We discuss two (related) Markov chains – the Gibbs sampler, and the Metropolis-Hastings algorithm.

107

## Markov chains over a discrete parameter space

Consider a random variable that may take on  $K$  values, and consider a Markov chain defined by a  $K \times K$  transition matrix  $\mathbf{P}$ .

Then the stationary distribution  $\pi$  is defined by

$$\pi = \pi \mathbf{P},$$

where  $\pi$  is a  $1 \times K$  row vector.

Roughly speaking, if  $\mathbf{P}$  is *irreducible* and *aperiodic* (i.e. ergodic) then the stationary distribution is unique.

## Gibbs Sampling

Consider a two-parameter problem in which the (intractable) posterior is:

$$\pi(\theta_1, \theta_2 | \mathbf{y}) \propto l(\theta_1, \theta_2) \times \pi(\theta_1, \theta_2).$$

We have

$$\pi(\theta_1, \theta_2 | \mathbf{y}) = p(\theta_1 | \mathbf{y}) \times p(\theta_2 | \theta_1, \mathbf{y}),$$

but  $p(\theta_1 | \mathbf{y})$  will typically be unavailable.

Gibbs sampling proceeds by iterating between the steps:

$$\theta_1^{(s)} \sim p(\theta_1 | \theta_2^{(s-1)}, \mathbf{y}),$$

and

$$\theta_2^{(s)} \sim p(\theta_2 | \theta_1^{(s)}, \mathbf{y}),$$

to produce the sequence

$$(\theta_1^{(0)}, \theta_2^{(0)}), (\theta_1^{(1)}, \theta_2^{(1)}), \dots, (\theta_1^{(s)}, \theta_2^{(s)}), \dots$$

which may be viewed as a draw from  $\pi(\theta_1, \theta_2 | \mathbf{y})$

109

## Gibbs Sampling over a discrete parameter space

Let  $\boldsymbol{\theta} = (\theta_1, \theta_2)$  and suppose that the parameters  $\theta_1$  and  $\theta_2$  can each take one of the two values, 0 and 1. The posterior distribution is given in Table 4.

$p(\theta_1, \theta_2   \mathbf{y})$	$\theta_2 = 0$	$\theta_2 = 1$
$\theta_1 = 0$	$\pi_{00}$	$\pi_{01}$
$\theta_1 = 1$	$\pi_{10}$	$\pi_{11}$

Table 4: Joint posterior distribution.

In this case the Gibbs sampler defines a  $4 \times 4$  transition matrix  $\mathbf{P}$ . The elements of this matrix are given by

$$\begin{aligned} \Pr\{(i, j), (k, l)\} &= \Pr\{\boldsymbol{\theta}^{(s)} = (k, l) | \boldsymbol{\theta}^{(s-1)} = (i, j)\} \\ &= \Pr(\theta_1^{(s)} = k | \theta_2^{(s)} = j) \Pr(\theta_2^{(s)} = l | \theta_1^{(s)} = k) \\ &= \frac{\pi_{kj}}{\pi_{+j}} \times \frac{\pi_{kl}}{\pi_{k+}} \end{aligned}$$

It is straightforward to show that  $\mathbf{P}$  is such that  $\pi = \pi \mathbf{P}$ .

### Example: Normal likelihood, unknown mean and variance

Likelihood:

$$Y_i | \boldsymbol{\beta}, \sigma^2 \sim N(\mathbf{x}_i \boldsymbol{\beta}, \sigma^2), i = 1, \dots, n.$$

Prior:

$$\boldsymbol{\beta} \sim N(\boldsymbol{\mu}, \mathbf{V}), \sigma^{-2} \sim \text{Ga}(a, b).$$

Posterior

$$\pi(\boldsymbol{\beta}, \sigma^2 | \mathbf{y}) \propto l(\boldsymbol{\beta}, \sigma^2) \pi(\boldsymbol{\beta}) \pi(\sigma^2),$$

is intractable unless  $p(\boldsymbol{\beta})$  is improper uniform and the prior for  $\sigma^2$  is inverse gamma.

111

Gibbs sampling iterates between  $\boldsymbol{\beta} | \mathbf{y}, \sigma^2$  and  $\sigma^{-2} | \mathbf{y}, \boldsymbol{\beta}$  where

$$\begin{aligned} p(\boldsymbol{\beta} | \mathbf{y}, \sigma^2) &\propto l(\boldsymbol{\beta}, \sigma^2) \pi(\boldsymbol{\beta}) \\ &\sim N(\boldsymbol{\mu}^*, \mathbf{V}^*), \\ p(\sigma^{-2} | \mathbf{y}, \boldsymbol{\beta}) &\propto l(\boldsymbol{\beta}, \sigma^2) \pi(\sigma^{-2}) \\ &\sim \text{Ga}\left(a + \frac{n}{2}, b + \frac{(\mathbf{y} - \mathbf{x}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{x}\boldsymbol{\beta})}{2}\right). \end{aligned}$$

where

$$\boldsymbol{\mu}^* = (\mathbf{x}^T \mathbf{x} \sigma^{-2} + \boldsymbol{\mu}^T \mathbf{V}^{-1})^{-1} (\mathbf{x}^T \mathbf{x} \hat{\boldsymbol{\beta}} \sigma^{-2} + \mathbf{b} \mathbf{V}^{-1}),$$

and

$$\mathbf{V}^* = (\mathbf{x}^T \mathbf{x} \sigma^{-2} + \mathbf{V}^{-1})^{-1}.$$

112

### Metropolis Algorithm – discrete parameter space

Suppose we have a discrete sample space  $\Omega$  and we wish to construct a Markov chain whose stationary distribution is  $\pi(\cdot)$ .

Let  $\mathbf{Q}$  be an irreducible transition matrix on  $\Omega$ , satisfying the symmetry condition

$$\mathbf{Q}(x, y) = \mathbf{Q}(y, x), \quad x, y \in \Omega.$$

We may then define a Markov chain  $\{\theta^{(s)}, s = 0, 1, 2, \dots\}$  via the following steps.

- Suppose we are currently at state  $x$ .
- Generate a proposal from  $\mathbf{Q}(x, y)$ .
- Accept  $\theta^{(s+1)} = y$  with probability

$$\min\left(1, \frac{\pi(y)}{\pi(x)}\right),$$

otherwise stay at  $x$ .

This results in the transition matrix

$$\mathbf{P}(x, y) = \mathbf{Q}(x, y) \times \min\left(1, \frac{\pi(y)}{\pi(x)}\right).$$

113

### Metropolis Algorithm – continuous parameter space

Suppose the stationary distribution is  $\pi(\boldsymbol{\theta})$  and consider the *symmetric* probability density function

$$g(\boldsymbol{\theta}_a | \boldsymbol{\theta}_b) = g(\boldsymbol{\theta}_b | \boldsymbol{\theta}_a).$$

Suppose  $\boldsymbol{\theta}^{(0)}$  denotes the initial point. The Metropolis algorithm then consists of, at iteration  $s$

- Sample  $\boldsymbol{\theta}^* | \boldsymbol{\theta}^{(s-1)} \sim g(\cdot | \boldsymbol{\theta}^{(s-1)})$ .
- Calculate  $r = \pi(\boldsymbol{\theta}^*) / \pi(\boldsymbol{\theta}^{(s-1)})$ .
- Set

$$\boldsymbol{\theta}^{(s)} = \begin{cases} \boldsymbol{\theta}^* & \text{with probability } \min(r, 1), \\ \boldsymbol{\theta}^{(s-1)} & \text{otherwise.} \end{cases}$$

At iteration  $s$  the transition density  $P(\boldsymbol{\theta}^{(s)} | \boldsymbol{\theta}^{(s-1)})$  is a mixture of  $g(\cdot | \boldsymbol{\theta}^{(s-1)})$  and the point  $\boldsymbol{\theta}^{(s-1)}$ .

Important point: the calculation of  $r$  does not depend on the normalizing constant of the target density  $\pi$ .

114

## Metropolis-Hastings Algorithm

Generalizes the Metropolis algorithm to allow a non-symmetric proposal density.

Suppose  $\boldsymbol{\theta}^{(0)}$  denotes the initial point. The Metropolis-Hastings algorithm then consists of, at iteration  $s$ :

- Sample  $\boldsymbol{\theta}^* | \boldsymbol{\theta}^{(s-1)} \sim g(\cdot | \boldsymbol{\theta}^{(s-1)})$ .
- Calculate

$$r = \frac{\pi(\boldsymbol{\theta}^*)/g(\boldsymbol{\theta}^* | \boldsymbol{\theta}^{(s-1)})}{\pi(\boldsymbol{\theta}^{(s-1)})/g(\boldsymbol{\theta}^{(s-1)} | \boldsymbol{\theta}^*)}.$$

- Set

$$\boldsymbol{\theta}^{(s)} = \begin{cases} \boldsymbol{\theta}^* & \text{with probability } \min(r, 1), \\ \boldsymbol{\theta}^{(s-1)} & \text{otherwise.} \end{cases}$$

115

## Issues:

- Convergence of the Markov chain?
- Parameterization.

## Convergence

- Early iterations  $\boldsymbol{\theta}^{(1)}, \boldsymbol{\theta}^{(2)}, \dots, \boldsymbol{\theta}^{(m)}$  reflect the (arbitrary) starting value  $\boldsymbol{\theta}^{(0)}$ .
- These iterations are called the *burn-in*.
- Chain will gradually ‘forget’ its initial state and converge to the unique stationary distribution which is independent of  $\boldsymbol{\theta}^{(0)}$ .
- Burn-in samples should be ignored when summarizing the samples for posterior inference via Monte Carlo integration, i.e.

$$E[g(\boldsymbol{\theta})] \approx \frac{1}{n-m} \sum_{s=m+1}^n g(\boldsymbol{\theta}^{(s)})$$

116



## Convergence Diagnosis

- Strictly speaking, convergence is only achieved for  $n = \infty$ .
- But we only need Markov chain to be ‘approaching’ convergence for Monte Carlo integration to yield a consistent estimate of the true expectation.
- How do we determine  $m$ , the number of ‘burn-in’ iterations?
- Informal examination of time series plots and running of multiple chains is a must.
- Two issues: have we ‘found’ the posterior? Do we have enough samples to answer the inferential questions? Some chains may be very slow mixing (examination of autocorrelation is important).

117

## Parameterization

The Markov chain will display better mixing properties if the parameters are approximately independent in the posterior.

In an extreme case, if we have independence then

$$p(\theta_1, \dots, \theta_k | \mathbf{y}) = \prod_{i=1}^k p(\theta_i | \mathbf{y}),$$

and Gibbs sampling via the conditional distributions  $p(\theta_i | \mathbf{y})$ ,  $i = 1, \dots, k$ , is equivalent to direct sampling from the posterior.

In general it is better to sample ‘blocks’ of parameters that are approximately independent.

118

## Hyperpriors

Consider the LMEM

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

with  $\mathbf{b}_i \sim N_{q+1}(\mathbf{0}, \mathbf{D})$ , and  $\boldsymbol{\epsilon}_i \sim N_{n_i}(\mathbf{0}, \sigma_\epsilon^2 \mathbf{I}_{n_i})$ ,  $i = 1, \dots, m$ . A Bayesian analysis requires prior distributions on  $\boldsymbol{\beta}, \mathbf{D}, \sigma_\epsilon^2$ ; it is common to assume independent priors

$$\pi(\boldsymbol{\beta}, \mathbf{D}, \sigma_\epsilon^2) = \pi(\boldsymbol{\beta})\pi(\mathbf{D})\pi(\sigma_\epsilon^2).$$

For  $\boldsymbol{\beta}$  a multivariate normal distribution and for  $\sigma_\epsilon^2$  an inverse gamma distribution are often specified since they lead to conditional distributions of convenient form for Gibbs sampling, but other choices are possible.

If  $\mathbf{D}$  is a diagonal matrix with elements  $\sigma_k^2$ ,  $k = 0, 1, \dots, q$ , then a prior that leads to conjugate conditional distributions in a Gibbs sampling algorithm is

$$\pi(\sigma_0^2, \dots, \sigma_q^2) = \prod_{k=0}^q \text{IGa}(a_k, b_k),$$

where  $\text{IGa}(a_k, b_k)$  denotes the inverse gamma distribution with pre-specified parameters  $a_k, b_k$ ,  $k = 0, \dots, q$ .

119

## The Wishart Distribution

A prior for a non-diagonal  $\mathbf{D}$  is more troublesome; there are  $(q+2)(q+1)/2$  elements, with the restriction that the resultant matrix is positive definite.

The inverse Wishart distribution is the conjugate choice, and is the only distribution for which any great practical experience has been gained.

Suppose  $\mathbf{Z}_1, \dots, \mathbf{Z}_r \sim_{iid} N_p(\mathbf{0}, \mathbf{S})$ , with  $\mathbf{S}$  a non-singular variance-covariance matrix, and let

$$\mathbf{W} = \sum_{j=1}^r \mathbf{Z}_j \mathbf{Z}_j^T. \quad (27)$$

Then  $\mathbf{W}$  follows a Wishart distribution, denoted  $W_p(r, \mathbf{S})$ , and

$$p(\mathbf{w}) = c^{-1} |\mathbf{w}|^{(r-p-1)/2} \exp\left\{-\frac{1}{2}\text{tr}(\mathbf{w}\mathbf{S}^{-1})\right\}$$

where

$$c = 2^{rp/2} \Gamma_p(r/2) |\mathbf{S}|^{r/2}, \quad (28)$$

with

$$\Gamma_p(r/2) = \pi^{p(p-1)/4} \prod_{j=1}^p \Gamma((r+1-j)/2)$$

the generalized gamma function, and  $r \geq p$  for a proper density.

The mean is given by

$$E[\mathbf{W}] = r\mathbf{S}.$$

The Wishart distribution is a multivariate version of the gamma distribution.

Taking  $p = 1$  yields

$$p(w) = \frac{(2S)^{-r/2}}{\Gamma(r/2)} w^{r/2-1} \exp(-w/2S),$$

for  $w > 0$ , the gamma distribution  $\text{Ga}(r/2, 1/(2S))$ . Further, taking  $S = 1$  gives a  $\chi_r^2$  random variable, which is clear from (27).

121

### The Inverse Wishart Distribution

If  $\mathbf{W} \sim W_p(r, \mathbf{S})$ , the distribution of  $\mathbf{D} = \mathbf{W}^{-1}$  is known as the inverse Wishart distribution, and is given by

$$p(\mathbf{d}) = c^{-1} |\mathbf{d}|^{-(r+p+1)/2} \exp\left\{-\frac{1}{2}\text{tr}(\mathbf{d}^{-1}\mathbf{S})\right\}$$

where  $c$  is again given by (28). The mean is given by

$$E[\mathbf{D}] = \frac{\mathbf{S}^{-1}}{r - p - 1}$$

and is defined for  $r > p + 1$ . If  $p = 1$  we recover the inverse gamma distribution  $\text{IGa}(r/2, 1/2S)$  with  $E[D] = 1/[s(r - 2)]$  and  $\text{var}(D) = 1/[S^2(r - 2)(r - 4)]$  (so that small  $r$  gives a larger spread).

Thinking ahead to application in the LMEM if  $\mathbf{W} \sim W_{q+1}(r, \mathbf{R}^{-1})$ , then

$$E[\mathbf{W}] = r\mathbf{R}^{-1},$$

and

$$E[\mathbf{D}] = \mathbf{R}/(r - q - 1 - 1),$$

so that  $\mathbf{R}$ , may be scaled to be a prior estimate of  $\mathbf{D}$ , with  $r$  acting as a strength of belief in the prior.

122

Issues with the Wishart Prior

- A problem with the Wishart distribution is that it is deficient in second moment parameters since there is only a single degrees of freedom parameter  $r$ . So, for example, it is not possible to have differing levels of certainty in the tightness of the prior distribution for different elements of  $\mathbf{D}$ . With diagonal  $\mathbf{D}$  and independent inverse gamma priors we have a precision parameter for each variance.
- The form of the conditional distribution suggests that it may be better to err on the side of picking  $\mathbf{R}$  too small (if  $m$  small, prior always influential).
- Intuition: as if our prior data for the precision consists of observing  $r$  normal random variables with variance-covariance matrices  $\mathbf{R}$ .
- We need to take  $r \geq q + 1$  for a proper prior, with the flattest prior corresponding to  $r = q + 1$ . A proper prior is required to ensure propriety of the posterior distribution.
- Figure 14 displays samples from the Wishart distribution  $W_2\{20, (20\mathbf{S})^{-1}\}$  where  $\mathbf{S} = \begin{bmatrix} 0.4 & 0 \\ 0 & 1.0 \end{bmatrix}$ . The mean is  $E[\mathbf{W}] = \mathbf{S}^{-1} = \begin{bmatrix} 2.5 & 0 \\ 0 & 1.0 \end{bmatrix}$ .

123

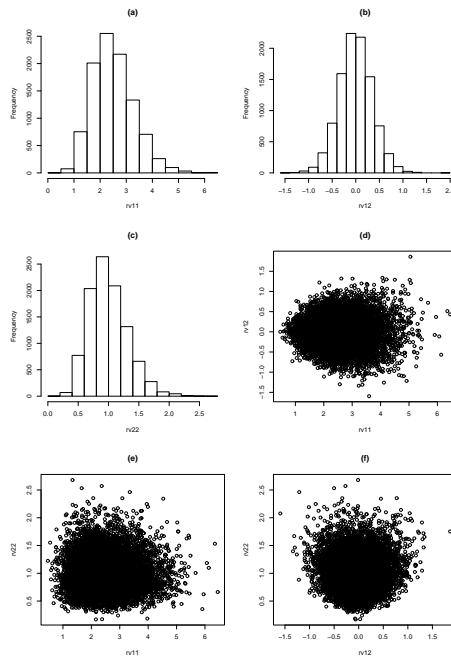


Figure 14: Histograms of (a)  $w_{11}$ , (b)  $w_{12}$ , (c)  $w_{22}$ , scatterplots of (d)  $w_{11}, w_{12}$ , (e)  $w_{11}, w_{22}$ ,  $w_{12}, w_{22}$

### Conditional Conjugacy

We now consider a Gibbs sampling scheme and assume for simplicity that  $\mathbf{x}_i = \mathbf{z}_i$ . It is computationally more convenient to reparameterize in terms of the set  $\{\boldsymbol{\beta}_1, \dots, \boldsymbol{\beta}_m, \tau, \boldsymbol{\beta}, \mathbf{W}\}$  where  $\boldsymbol{\beta}_i = \boldsymbol{\beta} + \mathbf{b}_i$ ,  $\tau = \sigma_\epsilon^{-2}$ ,  $\mathbf{W} = \mathbf{D}^{-1}$ .

The joint posterior is

$$p(\boldsymbol{\beta}_1, \dots, \boldsymbol{\beta}_m, \tau, \boldsymbol{\beta}, \mathbf{W}, \mathbf{b} \mid \mathbf{y}) \propto \prod_{i=1}^m \{p(\mathbf{y}_i \mid \boldsymbol{\beta}_i, \tau)p(\boldsymbol{\beta}_i \mid \boldsymbol{\beta}, \mathbf{W})\} \pi(\boldsymbol{\beta})\pi(\tau)\pi(\mathbf{W}),$$

with priors:

$$\begin{aligned} \boldsymbol{\beta} &\sim N_{q+1}(\boldsymbol{\beta}_0, \mathbf{V}_0) \\ \tau &\sim \text{Ga}(a_0, b_0) \\ \mathbf{W} &\sim W_{q+1}(r, \mathbf{R}^{-1}) \end{aligned}$$

and derive the required conditional distributions:

- $p(\boldsymbol{\beta} \mid \tau, \mathbf{W}, \boldsymbol{\beta}_1, \dots, \boldsymbol{\beta}_m, \mathbf{y})$
- $p(\tau \mid \boldsymbol{\beta}, \mathbf{W}, \boldsymbol{\beta}_1, \dots, \boldsymbol{\beta}_m, \mathbf{y})$
- $p(\mathbf{W} \mid \boldsymbol{\beta}, \tau, \boldsymbol{\beta}_1, \dots, \boldsymbol{\beta}_m, \mathbf{y})$
- $p(\boldsymbol{\beta}_i \mid \boldsymbol{\beta}, \tau, \mathbf{W}, \mathbf{y}), i = 1, \dots, m.$

125

### Conditional for $\boldsymbol{\beta}$

$$\boldsymbol{\beta} \mid \boldsymbol{\beta}_1, \dots, \boldsymbol{\beta}_m, \mathbf{W} \sim N_{q+1} \left\{ \left( m\mathbf{W} + \mathbf{V}_0^{-1} \right)^{-1} \left( \mathbf{W} \sum_{i=1}^m \boldsymbol{\beta}_i + \mathbf{V}_0^{-1} \boldsymbol{\beta}_0 \right), \left( m\mathbf{W} + \mathbf{V}_0^{-1} \right)^{-1} \right\}$$

### Conditional for $\tau$

$$\tau \mid \boldsymbol{\beta}_i, \mathbf{y} \sim \text{Ga} \left( a_0 + \frac{\sum_{i=1}^m n_i}{2}, b_0 + \frac{1}{2} \sum_{i=1}^m (\mathbf{y}_i - \mathbf{x}_i \boldsymbol{\beta}_i)^\top (\mathbf{y}_i - \mathbf{x}_i \boldsymbol{\beta}_i) \right)$$

### Conditional for $\boldsymbol{\beta}_i$

$$\boldsymbol{\beta}_i \mid \tau, \mathbf{W}, \mathbf{y} \sim N_{q+1} \left\{ (\tau \mathbf{x}_i^\top \mathbf{x}_i + \mathbf{W})^{-1} (\tau \mathbf{x}_i^\top \mathbf{y}_i + \mathbf{W} \boldsymbol{\beta}), (\tau \mathbf{x}_i^\top \mathbf{x}_i + \mathbf{W})^{-1} \right\}$$

Note the way that the conditional independencies have been exploited so that in each case we condition on only a subset of the parameters.

### Conditional for $\mathbf{W}$

First note that

$$(\boldsymbol{\beta}_i - \boldsymbol{\beta})^T \mathbf{W} (\boldsymbol{\beta}_i - \boldsymbol{\beta}) = \text{tr}((\boldsymbol{\beta}_i - \boldsymbol{\beta})^T \mathbf{W} (\boldsymbol{\beta}_i - \boldsymbol{\beta})) = \text{tr}(\mathbf{W} (\boldsymbol{\beta}_i - \boldsymbol{\beta}) (\boldsymbol{\beta}_i - \boldsymbol{\beta})^T).$$

Then

$$\begin{aligned} \mathbf{W} \mid \mathbf{y}, \boldsymbol{\beta}_i, \boldsymbol{\beta} &\propto \prod_{i=1}^m p(\boldsymbol{\beta}_i \mid \mathbf{W}) \times \pi(\mathbf{W}) \\ &\propto |\mathbf{W}|^{(m+r-q-1-1)/2} \exp \left\{ -\frac{1}{2} \left[ \sum_{i=1}^m (\boldsymbol{\beta}_i - \boldsymbol{\beta})^T \mathbf{W} (\boldsymbol{\beta}_i - \boldsymbol{\beta}) + \text{tr}(\mathbf{W} \mathbf{R}) \right] \right\} \\ &= |\mathbf{W}|^{(m+r-q-1-1)/2} \exp \left\{ -\frac{1}{2} \text{tr} \left( \mathbf{W} \left[ \sum_{i=1}^m (\boldsymbol{\beta}_i - \boldsymbol{\beta}) (\boldsymbol{\beta}_i - \boldsymbol{\beta})^T + \mathbf{R} \right] \right) \right\} \end{aligned}$$

Hence the conditional distribution is

$$\mathbf{W} \mid \boldsymbol{\beta}_1, \dots, \boldsymbol{\beta}_m, \boldsymbol{\beta}, \mathbf{y} \sim W_{q+1} \left\{ r + m, \left( \mathbf{R} + \sum_{i=1}^m (\boldsymbol{\beta}_i - \boldsymbol{\beta}) (\boldsymbol{\beta}_i - \boldsymbol{\beta})^T \right)^{-1} \right\}.$$

127

### Example: Dental Data for Girls

Three-Stage Hierarchical Model:

*First Stage:*

$$y_{ij} = \beta_{0i} + \beta_{1i}(t_j - 11) + \epsilon_{ij},$$

with  $\epsilon_{iid} \sim N(0, \tau^{-1})$ ,  $j = 1, \dots, 4$ ,  $i = 1, \dots, 11$ .

*Second Stage:* Let

$$\boldsymbol{\beta}_i = \begin{bmatrix} \beta_{0i} \\ \beta_{1i} \end{bmatrix} \quad \boldsymbol{\beta} = \begin{bmatrix} \beta_0 \\ \beta_1 \end{bmatrix} \quad \mathbf{D} = \begin{bmatrix} D_{00} & D_{01} \\ D_{10} & D_{11} \end{bmatrix},$$

and then

$$\boldsymbol{\beta}_i \mid \boldsymbol{\beta}, \mathbf{D} \sim N_2(\boldsymbol{\beta}, \mathbf{D}),$$

$i = 1, \dots, m$ .

*Third Stage:*

$$\pi(\tau, \boldsymbol{\beta}, \mathbf{D}^{-1}) \propto \text{Ga}(0, 0) \times N_2 \left( \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 10^6 & 0 \\ 0 & 10^6 \end{bmatrix} \right) \times W_2(r, \mathbf{R}^{-1}).$$

128

Results below are for priors, with prior mean

$$E[\mathbf{D}] = \frac{1}{r-q-2} \mathbf{R} = \frac{1}{r-3} \mathbf{R} = \begin{bmatrix} 1.0 & 0 \\ 0 & 0.1 \end{bmatrix}$$

(since  $q = 1$ ) and different degrees of freedom  $r$ .

We see sensitivity to the prior in inference for  $\mathbf{D}$ , but not for  $\boldsymbol{\beta}$ .

Note the greater shrinkage to the prior mean for the second and third priors.

$r$	$\mathbf{R}$	$\beta_0$	$\beta_1$
4	1.0 0 0 0.1	22.6 (21.4,23.8)	0.48 (0.33,0.63)
7	4.0 0 0 0.4	22.6 (21.5,23.7)	0.48 (0.31,0.65)
28	25 0 0 2.5	22.6 (21.8,23.5)	0.48 (0.28,0.67)

Table 5: Posterior medians and 95% intervals for population means, under three priors.

129

$r$	Diag $\mathbf{R}$	$D_{00}$	$D_{01}$	$D_{11}$
4	1.0 0.1	3.48 (1.66, 8.75)	0.13 (-0.10,0.54)	0.03 (0.01,0.10)
7	4.0 0.4	2.97 (1.51, 6.63)	0.10 (-0.14,0.46)	0.05 (0.02,0.12)
28	25 2.5	1.78 (1.14, 2.97)	0.04 (-0.10,0.20)	0.08 (0.05,0.14)

Table 6: Posterior medians and 95% intervals for population variances, under two priors.

130





## Quasi-Likelihood

We now describe a method for inference, generalized estimating equations, that attempts to make minimal assumptions about the data-generating process.

We begin with a recap of the related quasi-likelihood procedure, which is an alternative to MLE, when we do not wish to commit to specifying the full distribution of the data and we can assume independent data. The resultant estimators are known as quasi-MLE (QMLE).

The approach is based on specifying the first two moments of the data only, and assuming they take the form:

$$\begin{aligned} E[\mathbf{Y} \mid \boldsymbol{\beta}] &= \boldsymbol{\mu}(\boldsymbol{\beta}) \\ \text{cov}(\mathbf{Y} \mid \boldsymbol{\beta}) &= \alpha \mathbf{V}\{\boldsymbol{\mu}(\boldsymbol{\beta})\} \end{aligned}$$

where  $\boldsymbol{\mu}(\boldsymbol{\beta}) = [\mu_1(\boldsymbol{\beta}), \dots, \mu_n(\boldsymbol{\beta})]^\top$  represents the regression function and  $\mathbf{V}$  is a diagonal matrix (so the observations are uncorrelated), with

$$\text{var}(Y_i \mid \boldsymbol{\beta}) = \alpha V\{\mu_i(\boldsymbol{\beta})\},$$

and  $\alpha > 0$  a scalar which is independent of  $\boldsymbol{\beta}$ .

133

Consider the sum of squares

$$(\mathbf{Y} - \boldsymbol{\mu})^\top \mathbf{V}^{-1} (\mathbf{Y} - \boldsymbol{\mu}) / \alpha, \quad (29)$$

where  $\boldsymbol{\mu} = \boldsymbol{\mu}(\boldsymbol{\beta})$  and  $\mathbf{V} = \mathbf{V}(\boldsymbol{\beta})$ . To minimize this sum of squares there are two ways to proceed.

First approach: differentiate and obtain

$$-2\mathbf{D}^\top \mathbf{V}^{-1} (\mathbf{Y} - \boldsymbol{\mu}) / \alpha + (\mathbf{Y} - \boldsymbol{\mu})^\top \frac{\partial \mathbf{V}^{-1}}{\partial \boldsymbol{\beta}} (\mathbf{Y} - \boldsymbol{\mu}) / \alpha,$$

where  $\mathbf{D}$  is the  $n \times p$  matrix of derivatives with elements  $\partial \mu_i / \partial \beta_j$ ,  $i = 1, \dots, n; j = 1, \dots, p$ . Unfortunately the expectation of this expression is not zero, and so an inconsistent estimator of  $\boldsymbol{\beta}$  will result.

Second approach: pretend  $\mathbf{V}$  is not a function of  $\boldsymbol{\beta}$ , so that  $\hat{\boldsymbol{\beta}}$  is the root of:

$$\mathbf{D}(\hat{\boldsymbol{\beta}})^\top \mathbf{V}(\hat{\boldsymbol{\beta}})^{-1} \{\mathbf{Y} - \boldsymbol{\mu}(\hat{\boldsymbol{\beta}})\} / \alpha = \mathbf{0}.$$

134

As shorthand we write this estimating function as

$$\mathbf{U}(\boldsymbol{\beta}) = \mathbf{D}^T \mathbf{V}^{-1} \{\mathbf{Y} - \boldsymbol{\mu}\} / \alpha. \quad (30)$$

This estimating function is linear in the data and so its properties are straightforward to evaluate. In particular:

1.  $E[\mathbf{U}(\boldsymbol{\beta})] = \mathbf{0}$ .
2.  $\text{cov}\{\mathbf{U}(\boldsymbol{\beta})\} = \mathbf{D}^T \mathbf{V}^{-1} \mathbf{D} / \alpha$ .
3.  $-E \left[ \frac{\partial \mathbf{U}}{\partial \boldsymbol{\beta}} \right] = \text{cov}\{\mathbf{U}(\boldsymbol{\beta})\} = \mathbf{D}^T \mathbf{V}^{-1} \mathbf{D} / \alpha$ .

Applying the earlier result on properties of estimators arising from estimating functions:

$$(\mathbf{D}^T \mathbf{V}^{-1} \mathbf{D})^{1/2} (\hat{\boldsymbol{\beta}}_n - \boldsymbol{\beta}) \rightarrow_d N_p(\mathbf{0}, \alpha \mathbf{I}_p),$$

where we have so far assumed that  $\alpha$  is known.

Since the root of (30) does not depend on  $\alpha$ ,  $\hat{\boldsymbol{\beta}}$  is consistent regardless. For appropriate standard errors we require an estimator of  $\alpha$  however.

135

### Unknown $\alpha$

Since

$$E[(\mathbf{Y} - \boldsymbol{\mu})^T \mathbf{V}^{-1} (\boldsymbol{\mu}) (\mathbf{Y} - \boldsymbol{\mu})] = n\alpha,$$

an unbiased estimator of  $\alpha$  would be

$$\hat{\alpha} = (\mathbf{Y} - \boldsymbol{\mu})^T \mathbf{V}^{-1} (\boldsymbol{\mu}) (\mathbf{Y} - \boldsymbol{\mu}) / n,$$

a degrees of freedom corrected (but not in general, unbiased) estimate is given by the Pearson statistic divided by its degrees of freedom:

$$\hat{\alpha} = \frac{1}{n-p} \sum_{i=1}^n \frac{(Y_i - \hat{\mu}_i)^2}{V(\hat{\mu}_i)},$$

where  $\hat{\mu}_i = \hat{\mu}_i(\hat{\boldsymbol{\beta}})$ .

The asymptotic distribution that is used in practice is therefore given by

$$(\hat{\mathbf{D}}^T \hat{\mathbf{V}}^{-1} \hat{\mathbf{D}} / \hat{\alpha})^{1/2} (\hat{\boldsymbol{\beta}}_n - \boldsymbol{\beta}) \rightarrow_d N_p(\mathbf{0}, \mathbf{I}_p),$$

In general we may use sandwich estimation with quasi-likelihood. We have

$$\text{var}(\hat{\boldsymbol{\beta}}) = (\mathbf{D}^T \mathbf{V}^{-1} \mathbf{D})^{-1} \mathbf{D}^T \mathbf{V}^{-1} \text{var}(\mathbf{Y}) \mathbf{V}^{-1} \mathbf{D} (\mathbf{D}^T \mathbf{V}^{-1} \mathbf{D})^{-1} \alpha^2,$$

and  $\text{var}(\mathbf{Y})$  may be estimated by the diagonal matrix with elements  $(Y_i - \hat{\mu}_i)^2$ .

## Why “Quasi”

Integration of the quasi-score (30) gives

$$l(\mu, \alpha) = \int_y^\mu \frac{y-t}{\alpha V(t)} dt$$

which, if it exists, behaves like a log-likelihood. As an example, for the model  $E[Y] = \mu$  and  $\text{var}(Y) = \alpha\mu$  we have

$$l(\mu, \alpha) = \int_y^\mu \frac{y-t}{\alpha t} dt = \frac{1}{\alpha} [y \log \mu - \mu + c],$$

where  $c = -y \log y - y$  and  $y \log \mu - \mu$  is the log likelihood of a Poisson random variable.

The word “quasi” refers to the fact that the score may or not correspond to a probability function.

For example, the variance function  $\mu^2(1-\mu)^2$  does not correspond to a probability distribution.

137

## Example: Air Pollution Data

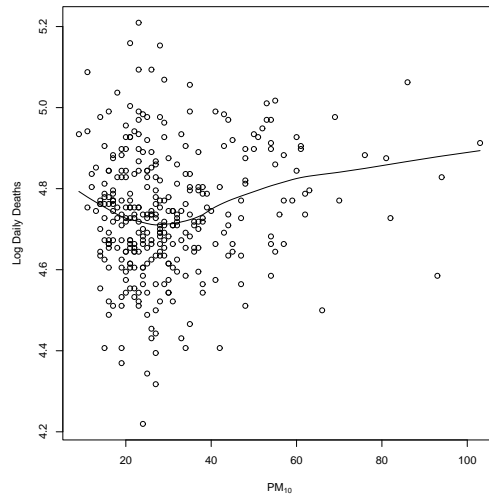
We examine the association between daily mortality,  $Y_i$ , and the daily value of  $\text{PM}_{10}$  (particulate matter less than 10 micrometers, which is about 0.0004 inches, in diameter),  $x_i$ , with  $i = 1, \dots, 335$ , indexing the 335 days on which there are no missing  $\text{PM}_{10}$  is to be investigated.

Figure 15 shows the association between log daily counts and  $\text{PM}_{10}$ .

Assume the model

$$E[Y_i | \boldsymbol{\beta}] = \exp(\mathbf{x}_i \boldsymbol{\beta}), \quad \text{var}(Y_i | \boldsymbol{\beta}) = \alpha E[Y_i | \boldsymbol{\beta}].$$

138

Figure 15: Log daily deaths versus  $PM_{10}$ .

139

Fitting the quasi-likelihood model yields  $\hat{\beta} = (4.71, 0.0015)^T$  and  $\hat{\alpha} = 2.77$  so that the quasi-likelihood standard errors are  $\sqrt{\hat{\alpha}} = 1.67$  times larger than the Poisson model-based standard errors.

The variance-covariance matrix is given by

$$(\hat{D}^T \hat{V}^{-1} \hat{D})^{-1} \hat{\alpha} = \begin{bmatrix} 0.019^2 & * \\ -0.89 \times 0.019 \times 0.00056 & 0.00056^2 \end{bmatrix}.$$

Standard errors of  $\hat{\beta}_0$  and  $\hat{\beta}_1$  are 0.019 and 0.00056.

Asymptotic 95% confidence interval for  $\beta_1$  is given by (0.00040, 0.0026).

A more useful summary is a confidence interval for the relative risk associated with a 10-unit increase in  $PM_{10}$ , which is

$$(e^{0.00040 \times 10}, e^{0.026 \times 10}) = (1.004, 1.026)$$

so that the interval suggests that the increase in daily mortality associated with a 10-unit increase in  $PM_{10}$  is between 0.4% and 2.6%.

140

## Extension to Quasi-Likelihood

Suppose we have

$$\begin{aligned} E[Y_i | \boldsymbol{\beta}] &= \mu_i(\boldsymbol{\beta}) \\ \text{var}(Y_i | \boldsymbol{\beta}) &= V_i(\boldsymbol{\alpha}, \boldsymbol{\beta}) \end{aligned}$$

where  $\boldsymbol{\alpha}$  is a  $k \times 1$  vector of parameters that appear only in the variance model.

Previously, in quasi-likelihood method, we had “separable” mean and variance models, that is,  $\text{var}(Y_i | \boldsymbol{\beta}) = \alpha V_i(\mu_i)$  (which is why we obtained a consistent estimator even if the form of the variance was wrong).

Let  $\hat{\boldsymbol{\alpha}}_n$  be a consistent estimator of  $\boldsymbol{\alpha}$ . We state without proof the following result. The estimator  $\hat{\boldsymbol{\beta}}_n$  that satisfies

$$\mathbf{G}(\hat{\boldsymbol{\beta}}_n, \hat{\boldsymbol{\alpha}}_n) = \mathbf{D}(\hat{\boldsymbol{\beta}}_n)^T \mathbf{V}^{-1}(\hat{\boldsymbol{\alpha}}_n, \hat{\boldsymbol{\beta}}_n) \{ \mathbf{Y} - \boldsymbol{\mu}(\hat{\boldsymbol{\beta}}_n) \} \quad (31)$$

has asymptotic distribution

$$(\hat{\mathbf{D}}^T \hat{\mathbf{V}}^{1/2} \hat{\mathbf{D}})^{-1} (\hat{\boldsymbol{\beta}}_n - \boldsymbol{\beta}) \rightarrow_d N_p(\mathbf{0}, \mathbf{I}_p) \quad (32)$$

where  $\hat{\mathbf{D}} = \mathbf{D}(\hat{\boldsymbol{\beta}}_n)$  and  $\hat{\mathbf{V}} = \mathbf{V}(\hat{\boldsymbol{\alpha}}_n, \hat{\boldsymbol{\beta}}_n)$ . Sandwich estimation may be used to obtain empirical standard errors which are correct even if the variance model is wrong, so long as we have a consistent estimator of  $\boldsymbol{\alpha}$ .

141

## Computation

Previously we assumed  $\text{var}(Y_i) = \alpha V_i(\mu_i)$ , and the estimating function did not depend on  $\alpha$  and so, correspondingly,  $\hat{\boldsymbol{\beta}}$  did not depend on  $\alpha$ , though the standard errors did.

In general iteration is convenient to simultaneously estimate  $\boldsymbol{\beta}$  and  $\boldsymbol{\alpha}$ .

Let  $\hat{\boldsymbol{\alpha}}^{(0)}$  be an initial estimate.

Then set  $j = 0$  and iterate between

1. Solve  $\mathbf{G}(\hat{\boldsymbol{\beta}}, \hat{\boldsymbol{\alpha}}^{(j)}) = \mathbf{0}$  to give  $\hat{\boldsymbol{\beta}}^{(j+1)}$ ,
2. Estimate  $\hat{\boldsymbol{\alpha}}^{(j+1)}$  with  $\hat{\mu}_i = \mu_i(\hat{\boldsymbol{\beta}}^{(j+1)})$ . Set  $j \rightarrow j + 1$  and return to 1.

142

### Example: Air Pollution Data

Consider the random effects formulation:

$$E[Y_i | \boldsymbol{\beta}, \theta_i] = \text{var}(Y_i | \boldsymbol{\beta}, \theta_i) = \mu_i(\boldsymbol{\beta})\theta_i \quad (33)$$

with

$$E[\theta_i] = 1, \quad \text{var}(\theta_i) = 1/\alpha. \quad (34)$$

Assuming  $\theta_i \sim_{iid} \text{Ga}(\alpha, \alpha)$ , we could derive the marginal distribution of the data (which is negative binomial) and proceed with likelihood.

As an alternative we consider the model

$$\begin{aligned} E[Y_i | \boldsymbol{\beta}] &= \mu_i(\boldsymbol{\beta}) \\ \text{var}(Y_i | \alpha, \boldsymbol{\beta}) &= \mu_i(\boldsymbol{\beta})\{1 + \mu_i(\boldsymbol{\beta})/\alpha\}. \end{aligned} \quad (35)$$

that are the marginal first two moments of the data given (33) and (34).

The form (35) suggests the estimating function for  $\boldsymbol{\beta}$  (with  $\alpha$  assumed known):

$$\sum_{i=1}^n \mathbf{D}(\boldsymbol{\beta})_i^T \mathbf{V}_i^{-1}(\alpha, \boldsymbol{\beta}) \{y_i - \mu_i(\boldsymbol{\beta})\}$$

For a fixed  $\alpha$  we can solve this estimating equation to obtain an estimator  $\widehat{\boldsymbol{\beta}}$ .

143

We describe a method-of-moments estimator for  $\alpha$  for the *quadratic* variance model we have

$$\text{var}(Y_i | \boldsymbol{\beta}, \alpha) = E[(Y_i - \mu_i)^2] = \mu_i(1 + \mu_i/\alpha),$$

and so

$$\alpha^{-1} = E \left[ \frac{(Y_i - \mu_i)^2 - \mu_i}{\mu_i^2} \right],$$

$i = 1, \dots, n$ , leading to the method-of-moments estimator

$$\widehat{\alpha} = \left\{ \frac{1}{n-p} \sum_{i=1}^n \frac{(Y_i - \widehat{\mu}_i)^2 - \widehat{\mu}_i}{\widehat{\mu}_i^2} \right\}^{-1}. \quad (36)$$

144

If we have a consistent estimator  $\hat{\alpha}$ , and the mean is correctly specified then valid inference follows from

$$(\hat{\mathbf{D}}^T \hat{\mathbf{V}}(\hat{\alpha})^{-1} \hat{\mathbf{D}})^{1/2} (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) \rightarrow_d N(\mathbf{0}, \mathbf{I}_p).$$

We fit this model to the air pollution data.

The estimates (standard errors) are  $\hat{\beta}_0 = 4.71$  (0.018) and  $\hat{\beta}_1 = 0.0014$  (0.00056).

The moment-based estimator is  $\hat{\alpha} = 65.20$ .

This analysis therefore produces virtually identical inference with the quasi-likelihood approach in which the variance was a linear function of the mean.

In Figure 16 we plot the linear and quadratic variance functions (over the range of the mean for these data) and we see that they are very similar.

Examination of the residuals did not clearly indicate the superiority of either variance model; it is typically very difficult to distinguish between the two models, unless the mean of the data has a large spread.

145

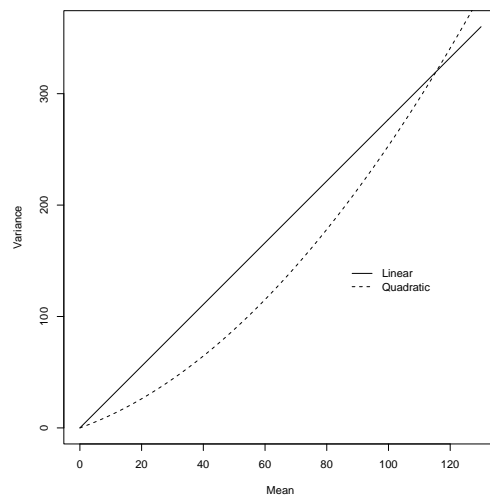


Figure 16: Linear and quadratic variance functions for the air pollution data.

146

### Example: Rcode for Quasi-Poisson Regression

We run the Poisson regression and then evaluate the method-of-moments estimator of  $\alpha$  “by hand”.

```
> mod1 <- glm(ynew~x1new,family=poisson)
> summary(mod1)
Coefficients:
                Value  Std. Error  t value
(Intercept) 4.705062304 0.0113962988 412.85880
      x1new 0.001458115 0.0003348748   4.35421
(Dispersion Parameter for Poisson family taken to be 1 )
      Null Deviance: 927.372 on 334 degrees of freedom
Residual Deviance: 908.6531 on 333 degrees of freedom
Number of Fisher Scoring Iterations: 3
Correlation of Coefficients:
      (Intercept)
x1new -0.8949913
> resid1 <- (ynew - mod1$fit)/sqrt(mod1$fit)
> alphahat <- sum(resid1 * resid1)/(length(ynew) - 2)
> alphahat
[1] 2.772861
```

147

We now fit the Quasi-Likelihood model with

$$E[Y_i|\beta] = \mu_i = \exp(\beta_0 + \beta_1 x_i)$$

and

$$\text{var}(Y_i|\beta) = \alpha\mu_i = \alpha \exp(\beta_0 + \beta_1 x_i).$$

```
> mod2 <- glm(ynew~x1new,quasi(link=log,variance=mu))
> summary(mod2)
Coefficients:
                Value  Std. Error  t value
(Intercept) 4.705062304 0.018976351 247.943468
      x1new 0.001458115 0.000557611   2.614932
(Dispersion Parameter for Quasi-likelihood
 family taken to be 2.772667 )
      Null Deviance: 927.372 on 334 degrees of freedom
Residual Deviance: 908.6531 on 333 degrees of freedom
Number of Fisher Scoring Iterations: 3
Correlation of Coefficients:
      (Intercept)
x1new -0.8949913
```

The standard errors are multiplied by  $\sqrt{\hat{\alpha}}$  (=1.67 here), but the estimates are unchanged.

148



### Example: Rcode for Quadratic Variance Model

The `glm.nb` function carries out MLE for the negative binomial model (it is part of the MASS library).

We find the MLE of  $\alpha$ , and then use this as a starting value for the iterative strategy in which a method-of-moments estimator is used.

```
> library(MASS)
> modnegbinmle <- glm.nb(y~x)
> summary(modnegbinmle)
Call:
glm.nb(formula = y ~ x, init.theta = 67.7145, link = log)
Coefficients:
            Estimate Std. Error z value Pr(>|z|)
(Intercept) 4.7055974  0.0188269 249.941  <2e-16 ***
x            0.0014405  0.0005577   2.583   0.0098 **
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
Correlation of Coefficients:
  (Intercept)
x -0.90
            Theta: 67.71
            Std. Err.: 8.27
> alphahat <- 67.61
```

149

Now iterate to a solution by estimating  $\beta$  for fixed  $\alpha$ , and then re-estimating  $\alpha$ .

```
> alphanew <- 0
> counter <- 0
> for (i in 1:5){
  fit <- glm(y~x,family=negative.binomial(alphahat))
  mu <- fit$fitted
  alphanew <- 1/(sum(((y-mu)^2-mu)/mu^2)/(length(y)-2))
  alphahat <- alphanew
  cat("Iteration ",i,alphahat,"\n")
}
Iteration 1 65.19642
Iteration 2 65.19649
Iteration 3 65.19649
Iteration 4 65.19649
Iteration 5 65.19649
> summary(fit)
Call:
glm(formula = y ~ x, family = negative.binomial(alphahat))
Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept) 4.705605  0.019071 246.747  <2e-16 ***
x            0.001440  0.000565   2.549   0.0112 *
(Dispersion parameter for Negative Binomial(65.1965)
family taken to be 1.001560)
```

150

## Generalized Estimating Equations

Suppose we assume

$$E[\mathbf{Y}_i | \boldsymbol{\beta}] = \mathbf{x}_i \boldsymbol{\beta},$$

and consider the  $n_i \times n_i$  *working* variance-covariance matrix:

$$\text{var}(\mathbf{Y}_i | \boldsymbol{\beta}, \boldsymbol{\alpha}) = \mathbf{W}_i.$$

To motivate GEE we begin by assuming that  $\mathbf{W}_i$  is known. In this case the GLS estimator minimizes

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

and is given by the solution to the estimating function

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

which is

$$\hat{\boldsymbol{\beta}} = \left( \sum_{i=1}^m \mathbf{x}_i^T \mathbf{W}_i^{-1} \mathbf{x}_i \right)^{-1} \sum_{i=1}^m \mathbf{x}_i^T \mathbf{W}_i^{-1} \mathbf{Y}_i.$$

We now examine the properties of this estimator.

151

We have

$$E[\hat{\boldsymbol{\beta}}] = \left( \sum_{i=1}^m \mathbf{x}_i^T \mathbf{W}_i^{-1} \mathbf{x}_i \right)^{-1} \sum_{i=1}^m \mathbf{x}_i^T \mathbf{W}_i^{-1} E[\mathbf{Y}_i] = \boldsymbol{\beta},$$

so long as the mean is correctly specified.

If the information about  $\boldsymbol{\beta}$  grows with increasing  $m$ , then  $\hat{\boldsymbol{\beta}}$  is consistent.

The variance,  $\text{var}(\hat{\boldsymbol{\beta}})$ , is given by

$$\left( \sum_{i=1}^m \mathbf{x}_i^T \mathbf{W}_i^{-1} \mathbf{x}_i \right)^{-1} \left( \sum_{i=1}^m \mathbf{x}_i^T \mathbf{W}_i^{-1} \text{var}(\mathbf{Y}_i) \mathbf{W}_i^{-1} \mathbf{x}_i \right) \left( \sum_{i=1}^m \mathbf{x}_i^T \mathbf{W}_i^{-1} \mathbf{x}_i \right)^{-1}.$$

If the assumed variance-covariance matrix is correct, i.e.  $\text{var}(\mathbf{Y}_i) = \mathbf{W}_i$ , then

$$\text{var}(\hat{\boldsymbol{\beta}}) = \left( \sum_{i=1}^m \mathbf{x}_i^T \mathbf{W}_i^{-1} \mathbf{x}_i \right)^{-1},$$

and a Gauss-Markov Theorem shows that, in this case, the estimator is efficient amongst linear estimators.

If  $m$  is large then a multivariate central limit theorem shows that  $\hat{\boldsymbol{\beta}}$  is asymptotically normal.

We now suppose that  $\text{var}(\mathbf{Y}_i) = \mathbf{W}_i(\boldsymbol{\alpha}, \boldsymbol{\beta})$  is of so that  $\boldsymbol{\alpha}$  are parameters in the variance-covariance model. The regression parameters are contained in  $\mathbf{W}_i$  to allow, mean-variance relationships, e.g.

$$\begin{aligned}\text{var}(Y_{ij} \mid \boldsymbol{\alpha}, \boldsymbol{\beta}) &= \alpha_1 \mu_{ij}^2 \\ \text{cov}(Y_{ij}, Y_{ik} \mid \boldsymbol{\alpha}, \boldsymbol{\beta}) &= \alpha_1 \alpha_2^{|t_{ij} - t_{ik}|} \mu_{ij} \mu_{ik}\end{aligned}$$

where

- $\mu_{ij} = \mathbf{x}_{ij} \boldsymbol{\beta}$ ,
- $\alpha_1$  is the variance (which is assumed constant across time and across individuals), and
- $\alpha_2$  is the correlation (which is assumed to be the same for all individuals), and
- $\boldsymbol{\alpha} = (\alpha_1, \alpha_2)$ .

153

For known  $\boldsymbol{\alpha}$  we would minimize

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

with solution given by the root of the estimating equation

$$\sum_{i=1}^m \mathbf{x}_i^T \mathbf{W}_i^{-1}(\boldsymbol{\alpha}, \boldsymbol{\beta}) (\mathbf{Y}_i - \mathbf{x}_i \boldsymbol{\beta}) = \mathbf{0}.$$

In general the roots of this equation are not available in closed form (because  $\boldsymbol{\beta}$  occurs in  $\mathbf{W}$ ).

However, if  $\mathbf{W}_i(\boldsymbol{\alpha}, \boldsymbol{\beta}) = \mathbf{W}_i(\boldsymbol{\alpha})$  we have

$$\hat{\boldsymbol{\beta}} = \left( \sum_{i=1}^m \mathbf{x}_i^T \mathbf{W}_i^{-1}(\boldsymbol{\alpha}) \mathbf{x}_i \right)^{-1} \sum_{i=1}^m \mathbf{x}_i^T \mathbf{W}_i^{-1}(\boldsymbol{\alpha}) \mathbf{Y}_i.$$

154

Finally, suppose that  $\boldsymbol{\alpha}$  is unknown but we have a method by which a consistent estimator  $\widehat{\boldsymbol{\alpha}}$  is produced (e.g. method of moments).

We then solve the estimator function

$$\mathbf{G}(\boldsymbol{\beta}) = \sum_{i=1}^m \mathbf{x}_i^T \mathbf{W}_i^{-1}(\widehat{\boldsymbol{\alpha}}, \boldsymbol{\beta})(\mathbf{Y}_i - \mathbf{x}_i \boldsymbol{\beta}).$$

In general iteration is needed to simultaneously estimate  $\boldsymbol{\beta}$  and  $\boldsymbol{\alpha}$ .

Let  $\widehat{\boldsymbol{\alpha}}^{(0)}$  be an initial estimate, then set  $t = 0$  and iterate between

1. Solve  $\mathbf{G}(\widehat{\boldsymbol{\beta}}, \widehat{\boldsymbol{\alpha}}^{(t)}) = \mathbf{0}$  to give  $\widehat{\boldsymbol{\beta}}^{(t+1)}$ ,
2. Estimate  $\widehat{\boldsymbol{\alpha}}^{(t+1)}$  with  $\widehat{\mu}_i = \mu_i(\widehat{\boldsymbol{\beta}}^{(t+1)})$ . Set  $t \rightarrow t + 1$  and return to 1.

155

We have

$$\text{var}(\widehat{\boldsymbol{\beta}})^{1/2}(\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta}) \sim N_{k+1}(\mathbf{0}, \mathbf{I}),$$

where

$$\begin{aligned} \widehat{\text{var}}(\widehat{\boldsymbol{\beta}}) &= \left( \sum_{i=1}^m \mathbf{x}_i^T \mathbf{W}_i^{-1}(\widehat{\boldsymbol{\alpha}}, \widehat{\boldsymbol{\beta}}) \mathbf{x}_i \right)^{-1} \\ &\times \left( \sum_{i=1}^m \mathbf{x}_i^T \mathbf{W}_i^{-1}(\widehat{\boldsymbol{\alpha}}, \widehat{\boldsymbol{\beta}}) \text{var}(\mathbf{Y}_i) \mathbf{W}_i^{-1}(\widehat{\boldsymbol{\alpha}}, \widehat{\boldsymbol{\beta}}) \mathbf{x}_i \right) \\ &\times \left( \sum_{i=1}^m \mathbf{x}_i^T \mathbf{W}_i^{-1}(\widehat{\boldsymbol{\alpha}}, \widehat{\boldsymbol{\beta}}) \mathbf{x}_i \right)^{-1}. \end{aligned}$$

We have assumed that  $\text{cov}(\mathbf{Y}_i, \mathbf{Y}_{i'}) = 0$  for  $i \neq i'$ , and this is required for the asymptotic distribution to be appropriate.

156

The final element of GEE is sandwich estimation of  $\text{var}(\widehat{\boldsymbol{\beta}})$ . In particular  $\text{cov}(\mathbf{Y}_i)$  is estimated by

$$(\mathbf{Y}_i - \mathbf{x}_i \widehat{\boldsymbol{\beta}})(\mathbf{Y}_i - \mathbf{x}_i \widehat{\boldsymbol{\beta}})^\top,$$

may be multiplied by  $N/(N-p)$  to account for estimation of  $\boldsymbol{\beta}$  ( $N = \sum_i n_i$ ).

*Empirical* would be a better word than *robust* (which is sometimes used) for the estimator of the variance – not robust to sample size, in fact could be highly unstable.

We can write the  $(k+1) \times 1$  estimating function as

$$\begin{aligned} & \mathbf{x}^\top \mathbf{W}^{-1}(\mathbf{Y} - \mathbf{x}\boldsymbol{\beta}) \\ & \sum_{i=1}^m \mathbf{x}_i^\top \mathbf{W}_i^{-1}(\mathbf{Y}_i - \mathbf{x}_i \boldsymbol{\beta}) \\ & \sum_{i=1}^m \sum_{j=1}^{n_i} [\mathbf{x}_{i1} \cdots \mathbf{x}_{in_i}] \begin{bmatrix} W_i^{11} & \cdots & W_i^{1n_i} \\ \cdots & \cdots & \cdots \\ W_i^{n_i 1} & \cdots & W_i^{n_i n_i} \end{bmatrix} \begin{bmatrix} Y_{i1} - \mathbf{x}_{i1} \boldsymbol{\beta} \\ \cdots \\ Y_{in_i} - \mathbf{x}_{in_i} \boldsymbol{\beta} \end{bmatrix} \end{aligned}$$

where  $W_i^{ij}$  denotes entry  $(i, j)$  of the inverse  $\mathbf{W}_i$ . We use the middle form since this emphasizes that the basic unit of replication is indexed by  $i$ .

157

**Example:** Suppose for simplicity that we have a balanced design, with  $n_i = n$  for all  $i$ , and assume a working variance-covariance matrix with

$$\begin{aligned} \text{var}(Y_{ij}) &= \text{E}[(Y_{ij} - \mathbf{x}_{ij} \boldsymbol{\beta})^2] = \text{E}[\epsilon_{ij}^2] = \alpha_1 \\ \text{cov}(Y_{ij}, Y_{ik}) &= \text{E}[(Y_{ij} - \mathbf{x}_{ij} \boldsymbol{\beta})(Y_{ik} - \mathbf{x}_{ik} \boldsymbol{\beta})] = \text{E}[\epsilon_{ij} \epsilon_{ik}] = \alpha_1 \alpha_{2jk}, \end{aligned}$$

for  $i = 1, \dots, m; j, k = 1, \dots, n; j \neq k$ . Hence we have  $n + n(n-1)/2$  elements of  $\boldsymbol{\alpha}$ .

Letting

$$e_{ij} = Y_{ij} - \mathbf{x}_{ij} \widehat{\boldsymbol{\beta}},$$

method-of-moments estimators are given by

$$\widehat{\alpha}_1 = \frac{1}{mn} \sum_{i=1}^m \sum_{j=1}^n e_{ij}^2,$$

and

$$\widehat{\alpha}_1 \widehat{\alpha}_{2jk} = \frac{1}{m} \sum_{i=1}^m e_{ij} e_{ik}.$$

158

*Generalized Estimating Equation (GEE) Summary*

We have:

- Regression parameters (of primary interest)  $\boldsymbol{\beta}$  and,
- Variance-covariance parameters  $\boldsymbol{\alpha}$ .

We have considered the GEE

$$\mathbf{G}(\boldsymbol{\beta}, \boldsymbol{\alpha}) = \sum_{i=1}^m \mathbf{D}_i^T \mathbf{W}_i^{-1} (\mathbf{Y}_i - \boldsymbol{\mu}_i) = \mathbf{0},$$

where

- $\boldsymbol{\mu}_i = \boldsymbol{\mu}_i(\boldsymbol{\beta}) = \mathbf{x}_i \boldsymbol{\beta}$ .
- $\mathbf{D}_i = \mathbf{D}_i(\boldsymbol{\beta}) = \frac{\partial \boldsymbol{\mu}_i}{\partial \boldsymbol{\beta}} = \mathbf{x}_i^T$ ,
- $\mathbf{W}_i = \mathbf{W}_i(\boldsymbol{\alpha}, \boldsymbol{\beta})$  is the “working” covariance model,

Three important ideas:

1. Separate estimation of  $\boldsymbol{\beta}$  and  $\boldsymbol{\alpha}$ .
2. Sandwich estimation of  $\text{var}(\hat{\boldsymbol{\beta}})$ .
3. Replication across units in order to estimate covariances – so we have assumed that observations on different units are independent.

159

Notes:

- We have seen the first and second ideas in independent data situations – e.g. estimation of the  $\alpha$  parameter in the quadratic negative binomial model.
- We may use method of moments estimators for  $\boldsymbol{\alpha}$  (or set up another estimating equation, see later).
- We could go with model-based standard errors:

$$\text{var}(\hat{\boldsymbol{\beta}}) = \left( \sum_{i=1}^m \mathbf{D}_i^T \mathbf{W}_i^{-1} \mathbf{D}_i \right)^{-1}. \quad (37)$$

If we have an independence working model ( $\mathbf{W}_i = \mathbf{I}$ ) then no iteration necessary (since no  $\boldsymbol{\alpha}$  in the GEE) – in this case we’d want to use sandwich estimation, however.

## Dental Example

Look at various estimators of  $\beta$  for girls only. Note here that we might question the asymptotics for GEE since we only have replication across  $m = 11$  units (girls) (check with simulation – see coursework).

Start with ordinary least squares – unbiased estimator for  $\beta$ , but standard errors are wrong because independence is assumed.

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

Call:

```
lm(formula = 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 **

Residual standard error: 2.164 on 42 degrees of freedom  
 Multiple R-Squared: 0.2046, Adjusted R-squared: 0.1856  
 F-statistic: 10.8 on 1 and 42 DF, p-value: 0.002053

161

Now implement GEE with working independence – the following is an R implementation.

```
> library(nlme); data(Orthodont); Orthgirl <- Orthodont[Orthodont$Sex=="Female",]
> install.packages("geepack")
> library(geepack)
> summary(geese(distance~age,id=Subject,data=Orthgirl,corstr="independence"))
Call:
geese(formula = distance ~ age, id = Subject, data = Orthgirl,corstr = "independence")
Mean Model:
  Mean Link:          identity
Variance to Mean Relation: gaussian
Coefficients:
      estimate   san.se     wald      p
(Intercept) 17.3727273 0.7819784 493.56737 0.000000e+00
age          0.4795455 0.0666386  51.78547 6.190604e-13
Scale Model:
  Scale Link:          identity
Estimated Scale Parameters:
      estimate   san.se     wald      p
(Intercept) 4.470403 1.373115 10.59936 0.001131270
Correlation Model:
  Correlation Structure:  independence
Returned Error Value:    0
Number of clusters:    11  Maximum cluster size: 4
```

162

Next we examine an exchangeable correlation structure in which all pairs of observations on the same unit have a common correlation:

```
> summary(geese(distance~age,id=Subject,data=Orthgirl,corstr="exchangeable"))
geese(formula = distance ~ age, id = Subject, data = Orthgirl,
      corstr = "exchangeable")
Mean Model:
Mean Link:          identity
Variance to Mean Relation: gaussian
Coefficients:
      estimate   san.se     wald      p
(Intercept) 17.3727273 0.7819784 493.56737 0.000000e+00
age          0.4795455 0.0666386  51.78547 6.190604e-13
Scale Model:
Scale Link:          identity
Estimated Scale Parameters:
      estimate   san.se     wald      p
(Intercept) 4.470403 1.373115 10.59936 0.001131270
Correlation Model:
Correlation Structure:  exchangeable
Correlation Link:      identity
Estimated Correlation Parameters:
      estimate   san.se     wald      p
alpha 0.8680178 0.1139327 58.04444 2.564615e-14
Number of clusters:    11  Maximum cluster size: 4
```

163

Notes:

- Independence estimates are always identical to OLS because we have assumed working independence, which means that the estimating equation is the same as the normal equations.
- Standard error for  $\beta_1$  is smaller with GEE because regressor (time) is changing within an individual.
- Here we obtain the same estimates for exchangeable as working independence but only because balanced and complete (i.e. no missing) data.



Finally we look at AR(1) and unstructured errors – this time we see slight differences in estimates and standard errors.

```
> summary(geese(distance~age,id=Subject,data=Orthgirl,corstr="ar1"))
geese(formula = distance ~ age, id = Subject, data = Orthgirl, corstr = "ar1")
Mean Model:
  Mean Link:          identity
  Variance to Mean Relation: gaussian
  Coefficients:
              estimate      san.se      wald      p
(Intercept) 17.3049830 0.85201953 412.51833 0.000000e+00
age          0.4848065 0.06881228  49.63692 1.849965e-12
Scale Model:
  Scale Link:          identity
  Estimated Scale Parameters:
              estimate      san.se      wald      p
(Intercept)  4.470639 1.341802 11.101 0.0008628115
Correlation Model:
  Correlation Structure:  ar1
  Correlation Link:      identity
  Estimated Correlation Parameters:
              estimate      san.se      wald p
alpha 0.9298023 0.07164198 168.4403 0
Number of clusters:  11  Maximum cluster size: 4
```

165

Now delete last two observations from girl 11 to illustrate that identical answers before were consequence of balance and completeness of data.

```
> Orthgirl2<-Orthgirl[1:42,]
> summary(lm(distance~age,data=Orthgirl2))
Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept)  18.0713      1.5102   11.966 8.56e-15 ***
age           0.3963      0.1357    2.921 0.00571 **
Residual standard error: 1.964 on 40 degrees of freedom
> summary(geese(distance~age,id=Subject,data=Orthgirl2,
corstr="independence"))
Coefficients:
              estimate      san.se      wald      p
(Intercept) 18.0713312 0.82603439 478.61250 0.000000e+00
age          0.3962971 0.06934195  32.66253 1.096304e-08
Scale Model:
  Scale Link:          identity
  Estimated Scale Parameters:
              estimate      san.se      wald      p
(Intercept)  3.674926 1.317669 7.778294 0.005287771
Correlation Model:
  Correlation Structure:  independence
Returned Error Value:  0
Number of clusters:  11  Maximum cluster size: 4
```

166

```

> summary(geese(distance~age,id=Subject,data=Orthgirl2,corstr="exchangeable"))
Call:
geese(formula = distance ~ age, id = Subject, data = Orthgirl2,
      corstr = "exchangeable")
Mean Model:
Mean Link:          identity
Variance to Mean Relation: gaussian
Coefficients:
              estimate      san.se      wald      p
(Intercept) 17.6050097 0.79007168 496.52320 0.000000e+00
age          0.4510122 0.06641218  46.11913 1.112765e-11
Scale Model:
Scale Link:          identity
Estimated Scale Parameters:
              estimate      san.se      wald      p
(Intercept) 3.706854 1.320019 7.88589 0.004982194
Correlation Model:
Correlation Structure:  exchangeable
Correlation Link:      identity
Estimated Correlation Parameters:
              estimate      san.se      wald p
alpha 0.7968515 0.09367467 72.36198 0
Returned Error Value:  0
Number of clusters:   11  Maximum cluster size: 4

```

167

### Comparison of Analyses

In Table 7 summaries are presented under likelihood, Bayesian and GEE analyses.

Two Bayesian models were fitted, a normal model:

$$\begin{aligned}
 \beta_i | \beta, \mathbf{D} &\sim_{iid} N(\beta, \mathbf{D}), \quad \text{var}(\beta_i | \beta, \mathbf{D}) = \mathbf{D} \\
 \mathbf{D}^{-1} &\sim W(r, \mathbf{R}^{-1}), \quad E[\text{var}(\beta_i | \beta, \mathbf{D})] = \frac{\mathbf{R}}{r-3} \\
 \mathbf{R} &= \begin{bmatrix} 1.0 & 0 \\ 0 & 0.1 \end{bmatrix}, \quad r = 4
 \end{aligned}$$

and a Student  $t_4$  model:

$$\begin{aligned}
 \beta_i | \beta, \mathbf{D} &\sim_{iid} \text{St}_4(\beta, \mathbf{D}), \quad \text{var}(\beta_i | \beta, \mathbf{D}) = 2\mathbf{D} \\
 \mathbf{D}^{-1} &\sim W(r, \mathbf{R}_t^{-1}), \quad E[\text{var}(\beta_i | \beta, \mathbf{D})] = 2\frac{\mathbf{R}_t}{r-3} \\
 \mathbf{R}_t &= \begin{bmatrix} 0.5 & 0 \\ 0 & 0.05 \end{bmatrix}, \quad r = 4
 \end{aligned}$$

168

Approach	$\widehat{\beta}_0$	s.e.( $\widehat{\beta}_0$ )	$\widehat{\beta}_1$	s.e.( $\widehat{\beta}_1$ )
LMEM ML	22.65	0.62	0.480	0.065
LMEM REML	22.65	0.63	0.479	0.066
Bayes Normal	22.65	0.60	0.479	0.075
Bayes $t_4$	22.65	0.58	0.475	0.073
GEE Independence	22.65	0.55	0.480	0.067
GEE AR(1)	22.64	0.58	0.485	0.069

Table 7: Summaries for fixed effects.

- Overall, the analyses are in good correspondence.

169

Approach	$\widehat{\text{var}}(\beta_{0i})$	$\widehat{\text{var}}(\beta_{1i})$	$\widehat{\text{corr}}(\beta_{0i}, \beta_{1i})$	$\widehat{\sigma}_\epsilon$
LMEM ML	1.98	0.15	0.55	0.67
LMEM REML	2.08	0.16	0.53	0.67
Bayes Normal	1.93 (1.29,2.96)	0.18 (0.10,0.31)	0.39 (-0.32,0.85)	0.70 (0.52,0.93)
Bayes $t_4$	2.06 (1.18,3.46)	0.20 (0.11,0.35)	0.42 (-0.34,0.88)	0.71 (0.54,0.95)

Table 8: Summaries for variance components.

GEE with working independence gives  $\alpha_1 = 4.47$ .

GEE with working AR(1) gives  $\alpha_1 = 4.47$ ,  $\alpha_2 = 0.93$ .

The parameterization adopted for the linear model changes the interpretation of  $\mathbf{D}$ . For example:

Model 1:  $(\beta_0 + b_{0i}) + (\beta_1 + b_{1i})t_j$ ,  $\mathbf{b}_i \sim \text{N}(\mathbf{0}, \mathbf{D})$ .

Model 2:  $(\gamma_0 + b_{0i}^*) + (\gamma_1 + b_{1i}^*)(t_j - \bar{t})$ ,  $\mathbf{b}_i^* \sim \text{N}(\mathbf{0}, \mathbf{D}^*)$ .

Giving  $\beta_0 = \gamma_0 - \gamma_1\bar{t}$ ,  $\beta_1 = \gamma_1$ .

$b_{0i} = b_{0i}^* - \bar{t}b_{1i}^*$ ,  $b_{1i} = b_{1i}^*$ .

Moral:  $\mathbf{D} \neq \mathbf{D}^*$ ;  $D_{00} = D_{00}^* - 2\bar{t}D_{01}^* + \bar{t}^2D_{11}^*$ ,  $D_{01} = D_{01}^* - \bar{t}D_{11}^*$ ,  $D_{11} = D_{11}^*$ .

170

### Covariance Models for Clustered Data

Whether we take a GEE or LME approach (with inference from the likelihood or from the posterior) we require flexible yet parsimonious covariance models.

In GEE we require a working covariance model

$$\text{cov}(\mathbf{Y}_i) = \mathbf{W}_i,$$

$i = 1, \dots, m.$

With LME we have so far assumed the model

$$\mathbf{y}_i = \mathbf{x}_i\boldsymbol{\beta} + \mathbf{z}_i\mathbf{b}_i + \boldsymbol{\epsilon}_i, \quad (38)$$

with  $\mathbf{b}_i \sim_{ind} N(\mathbf{0}, \mathbf{D})$  and  $\boldsymbol{\epsilon}_i \sim_{ind} N(\mathbf{0}, \mathbf{E}_i)$ , with  $\mathbf{E}_i = \mathbf{I}_{n_i}\sigma^2$ .

With  $\mathbf{z}_i\mathbf{b}_i = \mathbf{1}_{n_i}b_i$  we obtained an *exchangeable* (also known as compound symmetry):

$$\text{var}(\mathbf{Y}_i) = \sigma^2 \begin{bmatrix} 1 & \rho & \rho & \rho \\ \rho & 1 & \rho & \rho \\ \rho & \rho & 1 & \rho \\ \rho & \rho & \rho & 1 \end{bmatrix}$$

This model is particularly appropriate for clustered data with no time ordering (e.g. ANOVA).

171

An obvious extension for longitudinal data is to assume

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

with:

- Random effects  $\mathbf{b}_i \sim_{ind} N(\mathbf{0}, \mathbf{D})$ .
- Serial correlation  $\boldsymbol{\delta}_i \sim_{ind} N(\mathbf{0}, \mathbf{R}_i\sigma_\delta^2)$ , with  $\mathbf{R}_i$  an  $n_i \times n_i$  correlation matrix with elements

$$R_{ijj'} = \text{corr}(Y_{ij}, Y_{ij'} | \mathbf{b}_i),$$

$$j, j' = 1, \dots, n_i.$$

- Measurement error  $\boldsymbol{\epsilon}_i \sim_{ind} N(\mathbf{0}, \mathbf{I}_{n_i}\sigma_\epsilon^2)$ .

In general it is difficult to identify all three sources of variability – but the above provides a useful conceptual model.

See DHLZ, Chapter 5; Verbeke and Molenberghs, Chapter 10; Pinheiro and Bates, Chapter 5.

172

## Within-Unit Covariance Models

### Autoregressive errors

A widely-used time series model is the autoregressive, AR(1), process

$$\delta_{ij} = \rho\delta_{i,j-1} + u_{ij}, \quad (39)$$

for  $j \geq 2$ ,  $|\rho| \leq 1$  where  $u_{ij} \sim_{iid} N(0, \sigma_u^2)$  and are independent of  $\delta_{ik}$ ,  $k > 0$ . For LMEM we require a likelihood and hence the joint distribution of  $\boldsymbol{\delta}_i$ , for GEE the first two moments.

Repeated application of (39) gives, for  $k > 0$ ,

$$\delta_{ij} = u_{ij} + \rho u_{i,j-1} + \rho^2 u_{i,j-2} + \dots + \rho^{k-1} u_{i,j-k+1} + \rho^k \delta_{i,j-k}. \quad (40)$$

Assume the process has been running since  $j = -\infty$  and that it is ‘stable’ so that  $|\rho| < 1$  and the  $\delta_{ij}$  all have the same distribution.

Then, from (40)

$$\text{var}(\delta_{ij}) = \sigma_u^2(1 + \rho^2 + \rho^4 + \dots + \rho^{2(k-1)}) + \rho^{2k} \text{var}(\delta_{i,j-k}).$$

173

As  $k \rightarrow \infty$ , since  $\sum_{l=1}^{\infty} x^{l-1} = 1/(1-x)$ ,

$$\text{var}(\delta_{ij}) = \frac{\sigma_u^2}{(1-\rho^2)} = \sigma_\delta^2,$$

and, by substitution of (40),

$$\text{cov}(\delta_{ij}, \delta_{i,j-k}) = \text{E}[\delta_{ij}\delta_{i,j-k}] = \frac{\sigma_u^2 \rho^k}{(1-\rho^2)} = \sigma_\delta^2 \rho^k.$$

Hence under this model we have

$$\mathbf{R}_i = \begin{bmatrix} 1 & \rho & \rho^2 & \dots & \rho^{n_i-1} \\ \rho & 1 & \rho & \dots & \rho^{n_i-2} \\ \rho^2 & \rho & 1 & \dots & \rho^{n_i-3} \\ \dots & \dots & \dots & \dots & \dots \\ \rho^{n_i-1} & \rho^{n_i-2} & \rho^{n_i-3} & \dots & 1 \end{bmatrix}$$

as the correlation matrix for  $\boldsymbol{\delta}_i$ .

Often this model is written in the form

$$\text{cov}(Y_{ij}, Y_{ik}) = \sigma_\delta^2 \exp(-\phi d_{ijk}),$$

( $\rho = e^{-\phi}$ ) with  $d_{ijk} = |t_{ij} - t_{ik}|$  which is valid for unequally-spaced times also.

174

*Toeplitz*: Unstructured correlation:

$$\text{var}(\mathbf{Y}_i) = \sigma^2 \begin{bmatrix} 1 & \rho_1 & \rho_2 & \rho_3 \\ \rho_1 & 1 & \rho_1 & \rho_2 \\ \rho_2 & \rho_1 & 1 & \rho_1 \\ \rho_3 & \rho_2 & \rho_1 & 1 \end{bmatrix}$$

Heterogeneous versions with non-constant variance can also be fitted.

For example, the heterogeneous exchangeable model is given by:

$$\text{var}(\mathbf{Y}_i) = \begin{bmatrix} \sigma_1^2 & \rho\sigma_1\sigma_2 & \rho\sigma_1\sigma_3 & \rho\sigma_1\sigma_4 \\ \rho\sigma_2\sigma_1 & \sigma_2^2 & \rho\sigma_2\sigma_3 & \rho\sigma_2\sigma_4 \\ \rho\sigma_3\sigma_1 & \rho\sigma_3\sigma_2 & \sigma_3^2 & \rho\sigma_3\sigma_4 \\ \rho\sigma_4\sigma_1 & \rho\sigma_4\sigma_2 & \rho\sigma_4\sigma_3 & \sigma_4^2 \end{bmatrix}$$

Note that we should be careful when specifying the covariance structure – identifiability problems may arise if we try to be too flexible.

175

### Assessment of Assumptions

Each of the approaches to modeling that we have described depend upon assumptions concerning the structure of the data; to ensure that inference is appropriate we need to attempt to check that these assumptions are valid.

We first recap the assumptions:

*GEE*

Model:

$$\mathbf{Y}_i = \mathbf{x}_i\boldsymbol{\beta} + \mathbf{e}_i,$$

with working covariance model  $\text{var}(\mathbf{e}_i) = \mathbf{W}_i(\boldsymbol{\alpha})$ ,  $i = 1, \dots, m$ .

G1 Marginal model  $E[\mathbf{Y}_i] = \mathbf{x}_i\boldsymbol{\beta}$  is appropriate.

G2  $m$  is sufficiently large for asymptotic inference to be appropriate.

G3  $m$  is sufficiently large for robust estimation of standard errors.

G4 The working covariance  $\mathbf{W}_i(\boldsymbol{\alpha})$  is not far from the “true” covariance structure; if this is the case then the analysis will be very inefficient (standard errors will be much bigger than they need to be).

176

*LMEM via Likelihood Inference*

Model:

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

with  $\mathbf{b}_i \sim N(\mathbf{0}, \mathbf{D})$ ,  $\boldsymbol{\epsilon}_i \sim N(\mathbf{0}, \mathbf{E}_i)$ ,  $\mathbf{b}_i$  and  $\boldsymbol{\epsilon}_i$  independent ( $\mathbf{E}_i$  may have complex structure depending on both independent and dependent terms),  $i = 1, \dots, m$ .

L1 Mean model for fixed effects  $\mathbf{x}_i\boldsymbol{\beta}$  is appropriate.

L2 Mean model for random effects  $\mathbf{z}_i\mathbf{b}_i$  is appropriate.

L3 Variance model for  $\boldsymbol{\epsilon}_i$  is correct.

L4 Variance model for  $\mathbf{b}_i$  is correct.

L5 Normality of  $\boldsymbol{\epsilon}_i$ .

L6 Normality of  $\mathbf{b}_i$ .

L7  $m$  is sufficiently large for asymptotic inference to be appropriate.

*LMEM via Bayesian Inference*

Model as for LMEM, plus priors for  $\boldsymbol{\beta}$  and  $\boldsymbol{\alpha}$ .

Each of L1–L6 (asymptotic inference is not required if, for example, MCMC is used, though “appropriate” priors are needed).

177

**Overall strategy**

Before any formal modeling is carried out the data should be examined, in table and plot form, to see if the data have been correctly read in and to see if there are outliers.

For those individuals with sufficient data, individual-specific models should also be fitted, to allow examination of the appropriateness of initially hypothesized models in terms of the:

- linear component (which covariates, including transformations and interactions),
- and assumptions about the errors, such as constant variance and serial correlation.

Following fitting of marginal, mixed models, the assumptions should then be re-assessed, primarily through residual analysis.

## Residual Analysis

Residuals may be defined with respect to different levels of the model.

A vector of unstandardized *population-level* (marginal) residuals is given by

$$\mathbf{e}_i = \mathbf{Y}_i - \mathbf{x}_i \boldsymbol{\beta}.$$

A vector of unstandardized *unit-level* (Stage One) residuals is given by

$$\boldsymbol{\epsilon}_i = \mathbf{Y}_i - \mathbf{x}_i \boldsymbol{\beta} - \mathbf{z}_i \mathbf{b}_i.$$

The vector of random effects,  $\mathbf{b}_i$ , is also a form of (Stage Two) residual.

Estimated versions of these residuals are given by

$$\begin{aligned} \hat{\mathbf{e}}_i &= \mathbf{Y}_i - \mathbf{x}_i \hat{\boldsymbol{\beta}} \\ \hat{\boldsymbol{\epsilon}}_i &= \mathbf{Y}_i - \mathbf{x}_i \hat{\boldsymbol{\beta}} - \mathbf{z}_i \hat{\mathbf{b}}_i \end{aligned}$$

and  $\hat{\mathbf{b}}_i$ ,  $i = 1, \dots, m$ .

Recall from consideration of the ordinary linear model that estimated residuals have dependencies induced by the estimation procedure; in the dependent data context the situation is much worse as the “true” residuals have dependencies due to the dependent error terms of the models used.

Hence standardization is essential to remove the dependence.

179

### *Standardized Population Residuals*

If  $\mathbf{V}_i(\boldsymbol{\alpha})$  is the true error structure then

$$\text{var}(\mathbf{e}_i) = \mathbf{V}_i, \quad \text{and} \quad \text{var}(\hat{\mathbf{e}}_i) \approx \mathbf{V}_i(\hat{\boldsymbol{\alpha}}),$$

so that the residuals are dependent under the model, which means that it is not possible to check whether the covariance model is correctly specified (both form of the correlation structure and mean-variance model).

Plotting  $\hat{e}_{ij}$  versus  $x_{ij}$  may also be misleading due to the dependence within the residuals.

As an alternative, let  $\hat{\mathbf{V}}_i = \mathbf{L}_i \mathbf{L}_i^T$  be the Cholesky decomposition of  $\hat{\mathbf{V}}_i = \mathbf{V}_i(\hat{\boldsymbol{\alpha}})$ , the estimated variance-covariance matrix.



We can use this decomposition to form

$$\widehat{\mathbf{e}}_i^* = \mathbf{L}_i^{-1} \widehat{\mathbf{e}}_i = \mathbf{L}_i^{-1} (\mathbf{Y}_i - \mathbf{x}_i \widehat{\boldsymbol{\beta}}).$$

so that  $\text{var}(\mathbf{e}_i^*) \approx \mathbf{I}_{n_i}$ . We have the model

$$Y_i^* = \mathbf{x}_i^* \boldsymbol{\beta} + \mathbf{e}_i^*$$

where  $\mathbf{Y}_i^* = \mathbf{L}_i^{-1} \mathbf{Y}_i$ ,  $\mathbf{x}_i^* = \mathbf{L}_i^{-1} \mathbf{x}_i$ ,  $\mathbf{e}_i^* = \mathbf{L}_i^{-1} \mathbf{e}_i$ .

Hence plots of  $\widehat{e}_{ij}^*$  against columns of  $\mathbf{x}_{ij}^*$  should not show systematic patterns, *if* the assumed form is correct.

QQ plots of  $\widehat{e}_{ij}^*$  versus the expected residuals from a normal distribution can be used to assess normality (normal residuals are not required for GEE, but will help asymptotically).

Unstandardized versions will still be normally distributed if the  $\mathbf{e}_i$  are (since the  $e_{ij}^*$  are linear combinations of  $\mathbf{e}_i$ ), though the variances may be non-constant, and there may be strong dependence between different points.

The correctness of the mean-variance relationship can be assessed via examination of  $e_{ij}^{*2}$  versus  $\widehat{\mu}_{ij}^* = \mathbf{x}_{ij}^* \widehat{\boldsymbol{\beta}}$ .

Local smoothers can be added to plots to aid interpretation. Plotting symbols also useful – unit number, or observation number.

181

### Stage One Residuals

If  $\boldsymbol{\epsilon}_i \sim N(\mathbf{0}, \sigma_i^2 \mathbf{I}_{n_i})$  then residuals

$$\widehat{\boldsymbol{\epsilon}}_i = \mathbf{Y}_i - \mathbf{x}_i \widehat{\boldsymbol{\beta}} - \mathbf{z}_i \widehat{\mathbf{b}}_i$$

may be formed. Standardized versions are given by  $\widehat{\boldsymbol{\epsilon}}_i / \widehat{\sigma}_i$ .

The standardized versions should be used if the  $\sigma_i$  are unequal across  $i$ . Some uses:

- Plot residuals against covariates. Departures may suggest adding in covariates, both to  $\mathbf{x}_i$  and  $\mathbf{z}_i$ .
- To provide QQ plots – mean-variance relationship is more important to detect than lack of normality (so long as sample size is not small).
- assess constant variance assumption – one useful plot is versus  $\widehat{\mu}_{ij} = \mathbf{x}_{ij} \widehat{\boldsymbol{\beta}} + \mathbf{z}_{ij} \widehat{\mathbf{b}}_i$ .
- assess if serial correlation present in residuals

may be plotted against covariates to assess the form of the model, with QQ plots assessing normality of the measurement errors.

If  $\boldsymbol{\epsilon}_i \sim N(\mathbf{0}, \sigma_\epsilon^2 \mathbf{R}_i)$  with  $\mathbf{R}_i$  a correlation matrix then the residuals should be standardized, as with population residuals.

182

### Stage Two Residuals

Predictions of the random effects  $\hat{\mathbf{b}}_i$  may be used to assess assumptions associated with the random effects distribution, in particular:

- Are the random effects normally distributed?
- If we have assumed independence between random effects, does this appear reasonable?
- Is the variance of the random effects independent of covariates  $\mathbf{x}_i$ ?

It should be born in mind that interpretation of random effects predictions is more difficult since they are functions of the data.

Recall that  $\hat{\mathbf{b}}_i$  are shrinkage estimators, and hence assumptions about  $\mathbf{b}_i$  may not be reflected in  $\hat{\mathbf{b}}_i$ .

We may fit curves for particular individuals with  $n_i$  large, and then check the assumptions from these.

For the LMEM it is better to examine first and second stage residuals – population residuals are a mixture so if something wrong not clear at which stage there is trouble.

183

### Model refinement

A big problem is over-fitting in which models become too dataset-specific as they are refined on the basis of the examination of diagnostics.

In practice, if refinement is carried out through the fitting of alternative models (e.g. transformation of covariates, choice of distribution for the responses), then interval estimates will often be too narrow since they are produced by conditioning on the final model, and hence do not reflect the mechanism by which the model was selected.

From a frequentist standpoint estimators and test statistics should be examined via their long-run behaviour *given* the model-fitting process, including refinement. To be more explicit, let  $P$  denote the procedure by which a final model  $M$  is decided upon. Then suppose it is of interest to examine the bias of a statistic  $T$ ,

$$E[T|P] = E_{M|P}\{E[T|M]\}. \quad (41)$$

In general it will be incorrect to report  $E[T | \hat{M}]$  where  $\hat{M}$  is the final model chosen, since this does not reflect the procedure by which  $\hat{M}$  was chosen, but rather acts as if the final model is the “truth”.

184

From a Bayesian standpoint the same problem exists because the posterior distribution should reflect all sources of uncertainty and *a priori* all possible models that may be entertained should be explicitly stated, with prior distributions being placed upon different likelihoods and the parameters of these likelihoods; model averaging should then be carried out across the different possibilities.

One solution to this difficulty is to never refine the model for a given data set. This approach is operationally pure but pragmatically dubious (unless one is in the context of a randomized experiment) since we may obtain appropriate inference for a model that is a very poor description of the phenomenon under study.

The philosophy suggested here is to think as carefully as possible about the initial model class before the analysis proceeds, but after fitting to carry out model checking and refine the model in the face of *clear* model misspecification, with refinement ideally being carried out within distinct *a priori* known classes.

So that, for example, examining quantile-quantile plots for different  $t$  distributions and picking the one that produces the straightest line would not be a good idea. Inference then proceeds as if the final model were the one that were chosen initially. This is clearly a subjective procedure but can be informally justified via either philosophical approaches.

185

Under a frequentist approach inference follows from the behaviour of an estimator under repeated sampling from the true model, and if an initial model is clearly wrong on the basis of a residual plot (say), then it is very unlikely to be close to the “true” model and hence it is more appropriate to obtain properties of estimators under the assumed model. With reference to (41), if a model is chosen because it is clearly superior to the alternatives, then it may be reasonable to assume that  $E[T | P] \approx E[T | \hat{M}]$ , because  $\hat{M}$  would be consistently chosen in repeated sampling under these circumstances.

In a similar vein, under a Bayesian approach the above procedure is consistent with model-averaging but with the posterior model weight being concentrated upon the chosen model (since alternative models are only rejected on the basis of clear inadequacy). The aim is to provide probability statements, from either philosophical standpoints that are “honest” representations of uncertainty. The above approach is relevant to analyses that are more confirmatory in their outlook, as opposed to being used for prediction, or for more exploratory purposes (for example, to gain clues to models that may be appropriate for future data analyses).

If aim of analysis is simply exploratory, then we can do what we like (as soon as we quote CIs or significance levels ... trouble!).

186

Example: Dental Growth Curves – Initial Plots

We now present some initial plots for the dental data – should not be viewed as comprehensive.

- Initial plots: QQ plots of LS estimates, both univariate (Figure 17) and bivariate (Figure 18).
- Estimates of  $\sigma_\epsilon$ : 0.97, 0.59, 0.95, 0.30, 0.58, 0.43, 0.47, 0.19, 0.58, 0.85, 0.89 – not a great deal of variability, so common variance assumption seems reasonable.
- No apparent mean-variance relationship (Figure 19).
- Figure 20 shows that there are clear differences in intercepts, and some variability in slopes.

187

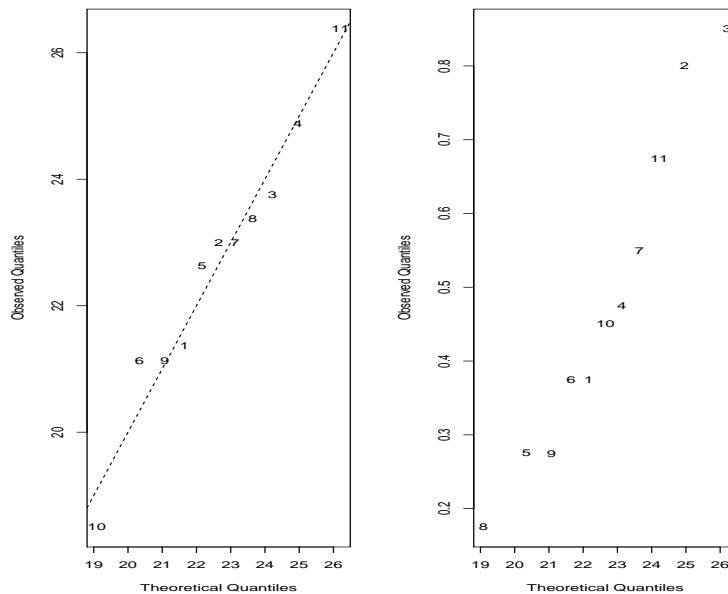


Figure 17: QQ plots of LS estimates:  $\hat{\beta}_0$  (left),  $\hat{\beta}_1$  (right).

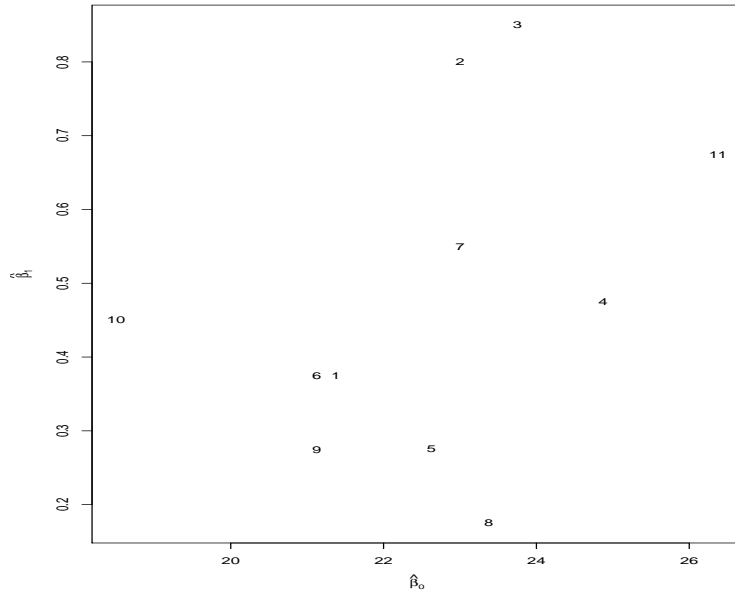


Figure 18: Bivariate plot of LS estimates.

189

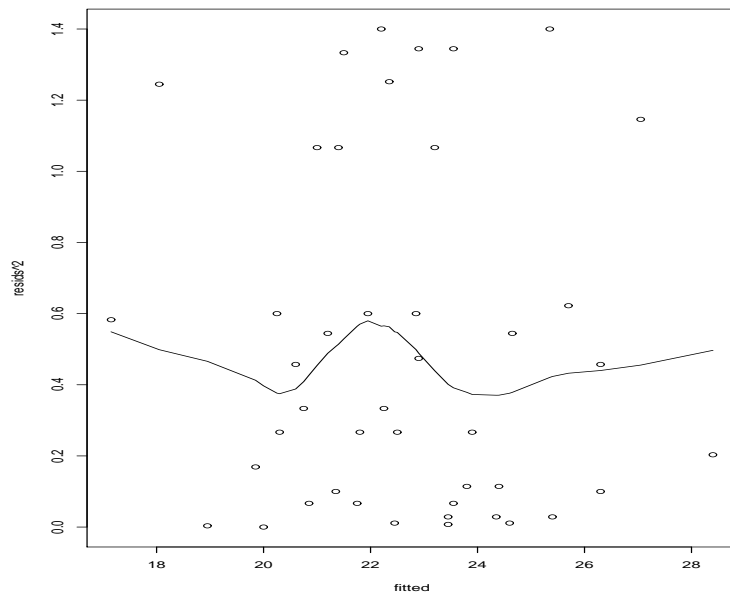


Figure 19: LS residuals versus fitted values

190

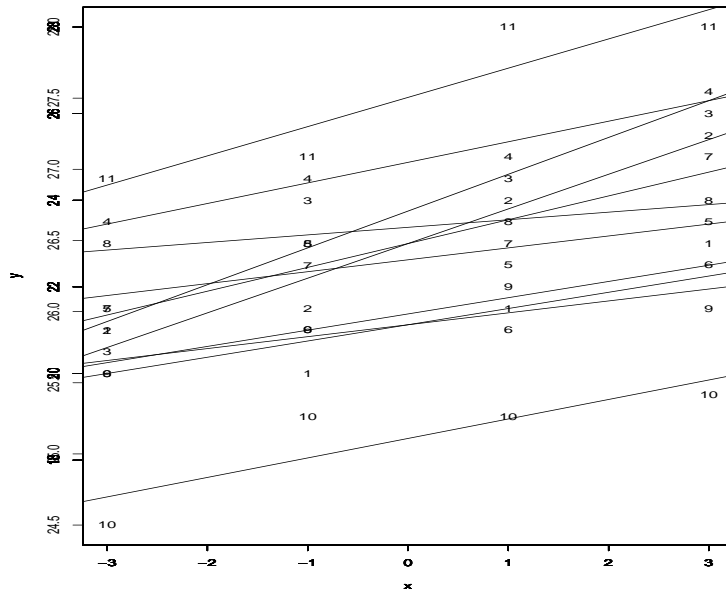


Figure 20: Fitted curves for all data.

191

### Example: Dental Growth Curves – Initial Plots

We now present some residual plots for the dental data.

```

> cnt1 <- rep(4,11); cnt2 <- 1:11
> lme1 <- lme(distance ~ I(age-11), data = Orthgirl, random = ~1 | Subject )
> lmeres1ind <- resid( lme1, level = 1, resType="n") # ind-level resids
> lmeFit1ind <- fitted( lme1 )
> plot(I(Orthgirl$age-11),lmeres1,xlab="Centered age",ylab="LME indiv residuals",
      ylim=c(-max(abs(range(lmeres1))),max(abs(range(lmeres1)))),type="n")
> text(I(Orthgirl$age-11),lmeres1,labels=c(rep(cnt2,cnt1))); abline(0,0)
> lines(lowess(I(Orthgirl$age-11),lmeres1))
> qqnorm(lmeres1,main="")
> plot(lmeFit1ind,lmeres1ind^2)
> lines(lowess(lmeFit1ind,lmeres1ind^2))

```

Figure 21 shows no systematic deviations between residuals with time. Figure 22 that normality reasonable, and Figure 23 that there is no mean variance relationship.

192

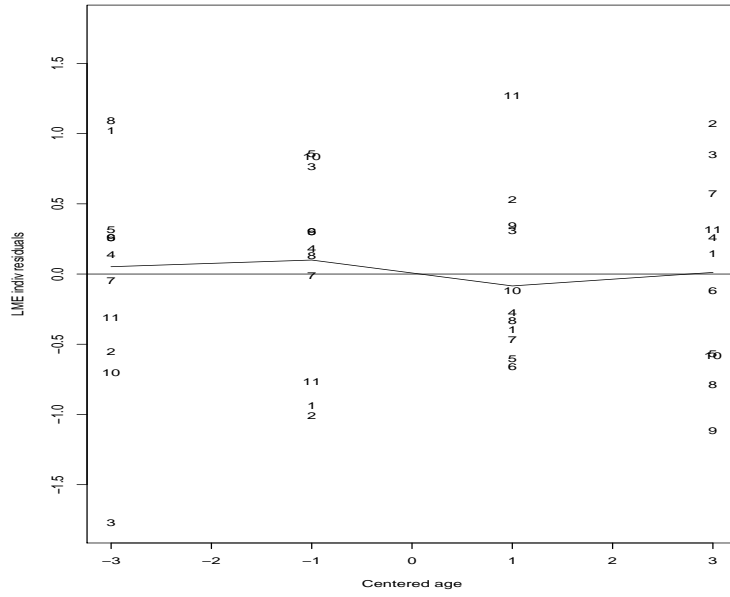


Figure 21: LME normalized residuals versus time.

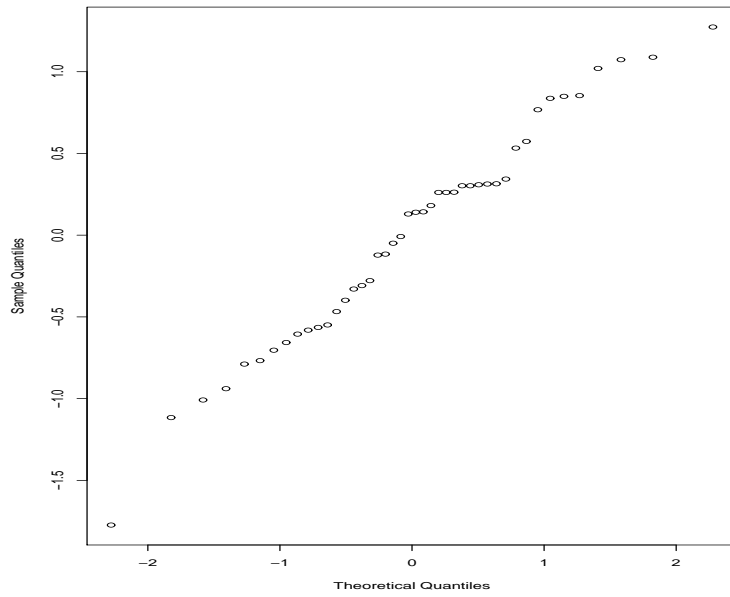


Figure 22: QQ plot of LME normalized residuals.

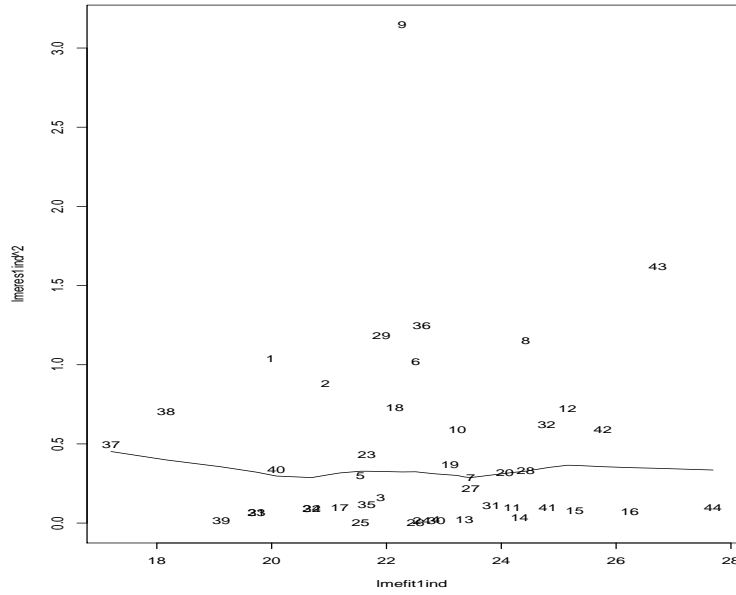


Figure 23: LME normalized residuals versus fitted values.

### Assessing Adequacy of the Temporal Covariance Structure

An informal method for assessing whether there is residual temporal dependence is to plot residuals versus time, we now consider more formal tools such as the correlogram and the variogram.

We begin with some definitions.

Consider a stochastic process  $Y(t)$  and let

$$\gamma(t, s) = \text{cov}\{Y(t), Y(s)\} = E[\{Y(t) - \mu(t)\}\{Y(s) - \mu(s)\}],$$

denote the *autocovariance function* of  $Y(t)$ .

The term *serial dependence* signifies that there is dependence between  $Y(t)$  and  $Y(s)$  for at least some pairs  $(s, t)$  with  $s \neq t$ .



We write

$$Y(t) = \mu(t) + e(t),$$

where  $\mu(t)$  is the deterministic trend component.

*Definition:* A process  $e(t)$  is second-order stationary if  $E[e(t)]$  is constant, for all  $t$ , and  $\gamma(t, s)$  depends only on  $|t - s|$ . For a residual process any non-zero constant has been absorbed into  $\mu(t)$ .

*Example:* The simplest example of a stationary random sequence is *white noise* which consists of a sequence of mutually independent random variables, each with mean 0 and finite variance  $\sigma^2$ .

There is a fundamental difficulty with trying to decompose  $Y(t)$  into the trend and the stochastic component in a single series because the two are unidentifiable without further assumptions.

Is it serial dependence in the residuals, or a high-order polynomial trend for example?

197

### The Autocorrelation Function

For a second-order stationary random process, the autocovariance function is

$$\text{cov}\{Y(t), Y(t + u)\} = \text{cov}\{e(t), e(t + u)\},$$

so that  $C(0)$  is the variance of  $Y(t)$  for all  $t$ .

The autocorrelation function is defined as

$$\rho(u) = \frac{C(u)}{C(0)}.$$

For equally-spaced data we could fit a model and then examine the autocorrelation function (ACF) of the residuals,

$$e_t = \frac{y_t - \hat{y}_t}{\widehat{\text{var}}(Y_t)^{1/2}}.$$

198

Consider a stochastic process  $e(t)$ , and realizations  $e_t$ ,  $t = 1, \dots, n$ . The empirical autocorrelation is defined as

$$\widehat{\rho}(u) = \widehat{\text{corr}}\{e(t), e(t+u)\} = \frac{\sum_{t=1}^{n-u} e_t e_{t+u} / (n-u)}{\sum_{t=1}^n e_t^2 / n},$$

for  $u = 0, 1, \dots$

A correlogram plot is  $\widehat{\rho}(u)$  versus  $u$ . If the residuals are a white noise process, we have the asymptotic result

$$\sqrt{n} e_t \rightarrow_d N(0, 1),$$

to give confidence bands  $\pm 2/\sqrt{n}$ .

199

## The Variogram

For unequally-spaced data the ACF is not so convenient, unless we round the observations.

An alternative is provided by the *semi-variogram* which is defined, for a process  $e_t$  and  $d \geq 0$ .

$$\gamma(d) = \frac{1}{2} \text{var}(e_t - e_{t-d}) = \frac{1}{2} \text{E} \left[ \{e_t - e_{t-d}\}^2 \right].$$

Recall that for a second-order stationary process,  $\text{E}[e_t] = \mu$  for all  $t$  and  $\text{cov}(e_t, e_{t-d})$  only depends on the distance  $d$  (which implies constant variance).

A smooth process is  $L_2$ -continuous, i.e.

$$\text{E}\{(e_t - e_{t-d})^2\} \rightarrow 0$$

as  $d \rightarrow 0$ . For a second-order stationary smooth process

$$\begin{aligned} \gamma(d) &= \frac{1}{2} \{ \text{E}[e_t^2] + \text{E}[e_{t-d}^2] - 2\text{E}[e_t e_{t-d}] \} \\ &= \sigma_e^2 \{1 - \rho(d)\}, \end{aligned}$$

where  $\text{var}(e) = \sigma_e^2$ .

200

The semi-variogram is also well-defined for an *intrinsically* stationary process for which  $E[e_t] = \mu$  and for which

$$E[(e_t - e_{t-d})^2] = 2\gamma(d).$$

As  $d$  increases then for observations far apart in time

$$\gamma(d) \rightarrow \text{var}(e_t) = \sigma_e^2,$$

which (recall) is assumed constant.

Consider measurement error,  $\epsilon_t$  with  $E[\epsilon_t] = 0$ ,  $\text{var}(\epsilon_t) = \sigma_\epsilon^2$ , and

$$Y_t = \mu_t + e_t + \epsilon_t,$$

so that we no longer have a smooth process. Then

$$\gamma(d) = \frac{1}{2}E[\{Y_t - Y_{t-d}\}^2] = \sigma_e^2\{1 - \rho(d)\} + \sigma_\epsilon^2,$$

and we have a “nugget” effect  $\sigma_\epsilon^2$ s.

201

### The Variogram in Longitudinal Data Analysis

Define the semi-variogram of the population residuals,  $e_{ij} = Y_{ij} - \mathbf{x}_{ij}\boldsymbol{\beta}$ , as

$$\gamma_i(d_{ijk}) = \frac{1}{2}E[\{e_{ij} - e_{ik}\}^2],$$

for  $d_{ijk} = |t_{ij} - t_{ik}| \geq 0$ . We emphasize that we are examining differences on the *same* individual.

The sample semi-variogram uses the empirical halved differences between pairs of population residuals

$$v_{ijk} = \frac{1}{2}(e_{ij} - e_{ik})^2,$$

along with the spacings  $u_{ijk} = t_{ij} - t_{ik}$ .

With highly-irregular sampling times the variogram can be estimated from the pairs  $(u_{ijk}, v_{ijk})$ ,  $i = 1, \dots, m$ ,  $j < k = 1, \dots, n_i$ , with the resultant plot being smoothed.

202

The marginal distribution of each  $v_{ijk}$  is  $\chi_1^2$ , and this large variability can make the variogram difficult to interpret.

The total variance is estimated as the average of  $\frac{1}{2}(e_{ij} - e_{lk})^2$ , for  $i \neq l$ , since

$$\frac{1}{2}\text{E}[(e_{ij} - e_{lk})^2] = \frac{1}{2}\{\text{E}[e_{ij}^2] + \text{E}[e_{lk}^2]\} = \sigma^2,$$

assuming that observations on different individuals are independent (and the variance is constant over time, and for different individuals).

Consider the interpretation of the variogram for the model

$$Y_{ij} = \mathbf{x}_{ij}\boldsymbol{\beta} + b_i + \delta_{ij} + \epsilon_{ij},$$

where  $b_i \sim_{ind} N(0, \sigma_0^2)$  (note, univariate),  $\epsilon_{ij} \sim_{ind} N(0, \sigma_\epsilon^2)$ , and  $\delta_{ij}$  represent error terms with serial dependence.

A simple and commonly-used form for serial dependence is the AR(1) model given by

$$\text{cov}(\delta_{ij}, \delta_{ik}) = \sigma_\delta^2 \rho^{|t_{ij} - t_{ik}|}.$$

Under this model

$$\text{var}(Y_{ij}|\boldsymbol{\beta}) = \sigma^2 = \sigma_0^2 + \sigma_\delta^2 + \sigma_\epsilon^2.$$

203

Consider the theoretical variogram for the residuals

$$e_{ij} = Y_{ij} - \mathbf{x}_{ij}\boldsymbol{\beta} = b_i + \delta_{ij} + \epsilon_{ij},$$

$i = 1, \dots, m; j = 1, \dots, n_i$ , with the AR(1) model.

For differences in residuals on the same individual

$$e_{ij} - e_{ik} = b_i + \delta_{ij} + \epsilon_{ij} - b_i - \delta_{ik} - \epsilon_{ik} = \delta_{ij} + \epsilon_{ij} - \delta_{ik} - \epsilon_{ik},$$

and so

$$\gamma_i(d_{ijk}) = \frac{1}{2}\text{E}[(e_{ij} - e_{ik})^2] = \sigma_\delta^2(1 - \rho^{d_{ijk}}) + \sigma_\epsilon^2. \quad (42)$$

As  $d_{ijk} \rightarrow 0$ ,  $\gamma_i(d_{ijk}) \rightarrow \sigma_\epsilon^2$  and  $b_i$  is the mean of  $e_{ij}$  and so its variance does not appear in (42).

Figure 24 shows the theoretical semi-variogram under this model and for the population residuals.

The variogram is limited in its use for *population* residuals for the LMEM, as we now illustrate.

Consider, the mixed effects model with random intercepts and independent random slopes:

$$b_{i0} \sim N(0, v_{00}^2), \quad b_{i1} \sim N(0, v_{11}^2)$$

leads to non-constant marginal variance

$$\text{var}(Y_{ij}|\boldsymbol{\beta}) = v_{00}^2 + 2v_{11}^2 t_{ij}^2,$$

so that we would not want to look at a variogram of population residuals because we do not have second-order stationarity. However, we could look at individual residuals after the random intercepts and slopes model has been fitted.

In my experience the variogram is often dominated by sampling variability (and there can be strong dependence in the plot since each residual contributes many points).

205

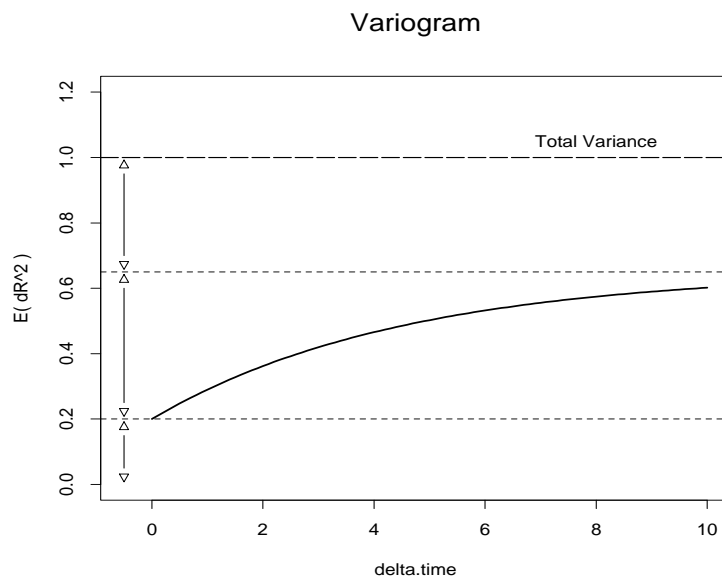


Figure 24: Theoretical variogram for a model with a random intercept, serial correlation, and measurement error.

206

Example: Air Pollution Data

We illustrate the correlogram and variogram for the air pollution data.

We fit a Poisson log-linear regression model in PM<sub>10</sub> and ozone.

In Figure 25 we clearly see strong dependence in the Pearson residuals, hence the quasi-likelihood standard errors quoted earlier will be wrong.

The dependence is confirmed by the dependence in the variogram in Figure 26. In the left-hand panel we have only plotted 1000 of the 53301 ( $327 \times 326/2$ ) points.

207

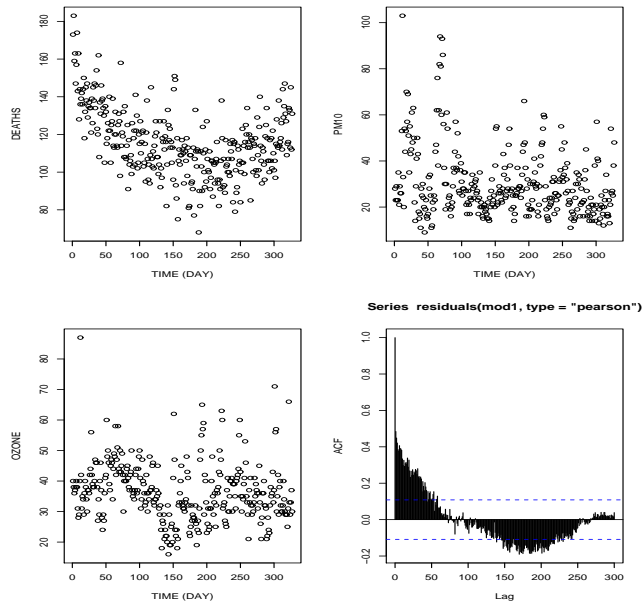


Figure 25: Time series plots and correlogram of residuals for air pollution data.

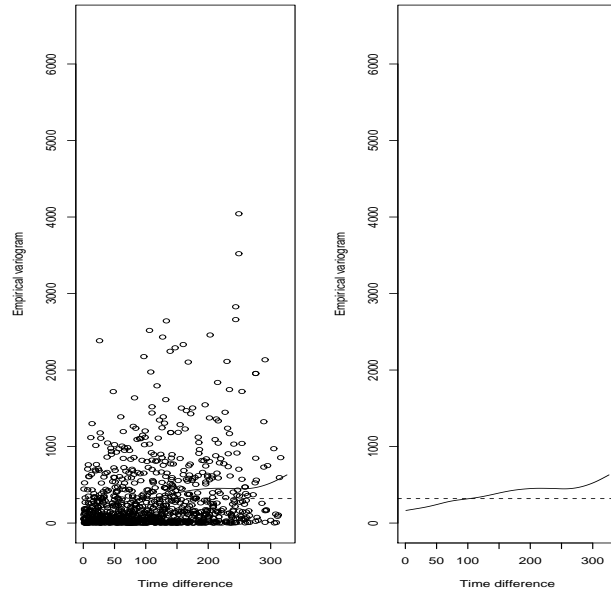


Figure 26: Variogram of residuals for air pollution data.