

# Statistical Methods for Correlated Data

## Likelihood Inference for Generalized Linear Mixed Models

Michele Guindani

Department of Biostatistics  
UCLA

- ▶ In GLMMs, there are three distinct sets of parameters for which inference may be required: fixed effects  $\beta$ , variance components  $\alpha$ , and random effects  $\mathbf{b} = [\mathbf{b}_1, \dots, \mathbf{b}_m]^T$ .
- ▶ **Objective:** maximize the likelihood  $L(\beta, \alpha)$ , where  $\alpha$  denote the variance components in  $\mathbf{D}$  and the scale parameter  $\phi$  (if present).
- ▶ The likelihood is obtained by integrating  $[\mathbf{b}_1, \dots, \mathbf{b}_m]$  from the model:

$$L(\beta, \alpha) = \prod_{i=1}^m \int p(\mathbf{y}_i | \beta, \mathbf{b}_i) \times p(\mathbf{b}_i | \alpha) d\mathbf{b}_i$$

- ▶ For non-Gaussian GLMMs, these integrals are not available in closed form.
- ▶ The resulting (approximated) likelihood may not be regular, so the second derivatives may be difficult to determine, so the Newton-Raphson method cannot be directly used  $\Rightarrow$  quasi-Newton approach (will discuss more later methods of estimating parameters)

# Procedures of maximization and hypothesis testing on fixed effects

- ▶ Given various estimating procedures for GLMMs,  $\beta$  is asymptotically normal, given the large-sample approximation to the joint distribution of parameters, with mean  $\mathbf{0}$  and variance  $\mathbf{I}_{\beta\beta}^{-1}$ .

- ▶ Briefly,

$$\begin{bmatrix} \hat{\beta} \\ \hat{\alpha} \end{bmatrix} \rightsquigarrow N \left( \begin{bmatrix} \beta \\ \alpha \end{bmatrix}, \begin{bmatrix} \mathbf{I}_{\beta\beta} & \mathbf{I}_{\beta\alpha} \\ \mathbf{I}_{\alpha\beta} & \mathbf{I}_{\alpha\alpha} \end{bmatrix}^{-1} \right)$$

- ▶ In general,  $\mathbf{I}_{\beta\alpha} \neq \mathbf{0}$  (compare with equation 8.20 in the textbook for LME) and so with GLMMs we cannot have REML estimates  
☞ consistency requires the specification of both mean and variance parameters

# Hypothesis testing

- ▶ If the asymptotic distribution is valid, we can use standard approaches for conducting hypothesis testing on the fixed effects.
- ▶ **Example, Wald test:** for testing  $H_0 : \hat{\beta} = \beta$  we can use

$$(\hat{\beta} - \beta)' \mathbf{I}(\hat{\beta})(\hat{\beta} - \beta) \sim \chi_p^2$$

where  $\mathbf{I}(\hat{\beta})$  is the asymptotically consistent estimator for the expected information matrix  $\mathbf{I}(\beta)$

# Hypothesis testing

- ▶ If the asymptotic distribution is valid, we can use standard approaches for conducting hypothesis testing on the fixed effects.
- ▶ **Example, Wald test:** for testing  $H_0 : \hat{\beta} = \beta$  we can use

$$(\hat{\beta} - \beta)' I(\hat{\beta})(\hat{\beta} - \beta) \sim \chi_p^2$$

where  $I(\hat{\beta})$  is the asymptotically consistent estimator for the expected information matrix  $I(\beta)$

- ▶ **Example, LRT:** can be used to test a null reduced model equation, as the likelihood ratio test statistics is asymptotically distributed as a  $\chi^2$  with degrees of freedom = to the difference in number of parameters in the model
- ▶ Between the likelihood ratio and the Wald statistics, the likelihood ratio statistic is considered to provide more reliable test results although it is computationally more complex. Therefore, when possible the likelihood ratio test is preferable to test a null hypothesis on the fixed effects.

# Hypothesis testing on the variance components

- ▶ When the significance test is performed on the variance-covariance components of the random effects, and the test is at the boundary of the parameter space, e.g.  $H_0 : \sigma_{b_i}^2 = 0$  and  $H_1 : \sigma_{b_i}^2 > 0$  where  $\sigma_{b_i}$  is the  $i$ -th diagonal element in  $\mathbf{D}$ , then the conventional tests, such as the likelihood ratio and the Wald statistics, cannot be directly applied to test these hypotheses.
- ▶ Since the  $\mathbf{D}$  matrix in GLMMs is generally specified to follow normality, one can still use the procedures described for LMEs and use mixtures of  $\chi^2$  distributions to test the null

# Methods for estimating parameters in GLMMs

- ▶ In the analysis of nonnormal longitudinal data, such as proportions or counts, numerical integration or Bayes-type techniques are required to conduct a full ML analysis.
- ▶ There are a variety of methods used to approximate the relevant likelihood quantities:
  1. Gaussian quadrature rules
  2. Laplace approximation
  3. Penalized quasi-likelihood (PQL)
  4. Marginal quasi-likelihood (MQL)
  5. Monte Carlo EM
  6. Markov Chain Monte Carlo (MCMC)

# Gaussian quadrature methods

- ▶ A widely used approximation method of the integrals in GLMMs is Gaussian quadrature
- ▶ Suppose that there is a known, smooth function  $f(z)$  and a probability density function  $p(z)$ . The function  $f(z)$  can then be integrated against  $p(z)$ . The corresponding quadrature rule is

$$\int_{-\infty}^{\infty} f(z)p(z)dz \approx \sum_{q=1}^Q w_q f(z_q)$$

where  $Q$  represents the number of quadrature points,  $w_q$  is the quadrature weight ( $q = 1, \dots, Q$ ), and  $z_q$  is the abscissas, statistically referred to as a node.



# Gaussian quadrature methods

- ▶ In GLMMs,  $p(z)$  is represented by the distribution of the random effects. However, in Gaussian quadratures,  $p(z)$  is standard normal, so the the random effects may first be standardized

$$\tilde{\mathbf{b}}_i = \mathbf{D}^{-1/2} \mathbf{b}_i$$

so that  $\tilde{\mathbf{b}} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ . Correspondingly, one needs to consider  $\boldsymbol{\eta}_i = \mathbf{X}_i' \boldsymbol{\beta} + \mathbf{Z}_i' \mathbf{D}^{1/2} \mathbf{b}_i$

- ▶ Given the above specification, the likelihood contribution for subject  $i$  is

$$\begin{aligned} f_i(\mathbf{y}_i | \boldsymbol{\beta}, \mathbf{D}, \phi) &= \int \prod_{j=1}^{n_i} f_{ij}(y_{ij} | \boldsymbol{\beta}, \mathbf{b}_i, \phi) f(\mathbf{b}_i | \mathbf{G}) d\mathbf{b}_i \\ &= \int \prod_{j=1}^{n_i} f_{ij}(y_{ij} | \boldsymbol{\beta}, \hat{\mathbf{b}}_i, \mathbf{D}, \phi) f(\hat{\mathbf{b}}_i) d\hat{\mathbf{b}}_i \end{aligned}$$

- ▶ A widely used algorithm is the Gauss-Hermite quadrature (details omitted)

# Laplace Approximation

- ▶ The Laplace method is a statistical technique for approximating integrals of the form

$$\int_{\hat{a}}^{\hat{b}} \exp[Nf(z)]dz$$

where  $f(z)$  is some known, smooth function that is unimodal with a maximum at  $z_0$  and  $N$  is the sample size.

- ▶ Therefore,  $z_0$  is assumed to be the only point satisfying

$$\Delta f(z) = \left[ \frac{\partial f(z)}{\partial z_{\hat{a}}}, \frac{\partial f(z)}{\partial z_{\hat{a}+1}}, \dots, \frac{\partial f(z)}{\partial z_{\hat{b}}} \right] = 0$$

- ▶ It is also assumed that the Hessian matrix of  $f(z)$  at  $z_0$

$$\mathbf{H}'_0 = \left[ \frac{\partial^2 f(z)}{\partial z_i \partial z_j} \right] \Big|_{z=z_0}$$

is positive definite

# Laplace Approximation

- ▶ According to Taylor's theorem,  $f(z)$  can be further expanded to the expression

$$f(z) = f(z_0) + f'(z_0)(z - z_0) + \frac{1}{2}f''(z_0)(z - z_0)^2 + O\left[(z - z_0)^3\right]$$

where  $O\left[(z - z_0)^3\right] \xrightarrow{p} 0$ .

- ▶ As  $f'(z_0) = 0$  by assumption,  $f(z)$  can be approximated to the quadratic form

$$f(z) = f(z_0) + \frac{1}{2}f''(z_0)(z - z_0)^2$$

- ▶ Then, assuming  $f''(z_0) < 0$ ,

$$\int_{\hat{a}}^{\hat{b}} \exp[Nf(z)]dz \approx \exp\{Nf(z_0)\} \int_{\hat{a}}^{\hat{b}} \exp\left\{-\frac{N}{2}\{-f''(z_0)\}(z - z_0)^2\right\} dz$$

which is a Normal kernel ( $N(z_0, \{-f''(z_0)\}^{-1}/N)$ ).

# Laplace approximation

► Hence,

$$\int_{\hat{a}}^{\hat{b}} \exp[N f(z)] dz \approx \exp\{N f(z_0)\} \sqrt{\frac{2\pi N}{\{-f''(z_0)\}}} \int_{\hat{a}}^{\hat{b}} \phi(z|z_0, \{-f''(z_0)\}^{-1}/N) dz$$

# Laplace approximation

- ▶ Hence,

$$\int_{\hat{a}}^{\hat{b}} \exp[N f(z)] dz \approx \exp\{N f(z_0)\} \sqrt{\frac{2\pi N}{\{-f''(z_0)\}}} \int_{\hat{a}}^{\hat{b}} \phi(z|z_0, \{-f''(z_0)\}^{-1}/N) dz$$

- ▶ The Laplace approx can be easily extended to the longitudinal setting.

$$\begin{aligned} p(\mathbf{y}) &= \prod_{i=1}^m p(\mathbf{y}_i) = \prod_{i=1}^m \int p(\mathbf{y}_i | \mathbf{b}_i) p(\mathbf{b}_i) d\mathbf{b}_i \\ &= \prod_{i=1}^m \int \exp\{\log[p(\mathbf{y}_i | \mathbf{b}_i) p(\mathbf{b}_i)]\} d\mathbf{b}_i \\ &= \prod_{i=1}^N \int \exp\left\{\sum_{j=1}^{n_i} \log[p(y_{ij} | \mathbf{b}_i)] + n_i \log[p(\mathbf{b}_i)]\right\} d\mathbf{b}_i \end{aligned}$$

# Laplace Approximation

- After some algebra, the Laplace approximation to the marginal log-likelihood can be written as

$$l(\boldsymbol{\beta}, \boldsymbol{\alpha}; \hat{\mathbf{b}}, \mathbf{y}) = \sum_{i=1}^N \left[ \sum_{j=1}^{n_i} \log \left[ p \left( y_{ij} | \hat{\mathbf{b}}_i \right) \right] + n_i \log \left[ p \left( \hat{\mathbf{b}}_i \right) \right] \right. \\ \left. + \frac{1}{2} n_{b_i} \log(2\pi) - \log \left| -\frac{1}{2} n_i f'' \left( \boldsymbol{\beta}, \boldsymbol{\alpha}; \hat{\mathbf{b}}_i \right) \right| \right]$$

where  $n_{b_i}$  is the dimension of the random effect  $\mathbf{b}_i$

# Laplace Approximation

- ▶ After some algebra, the Laplace approximation to the marginal log-likelihood can be written as

$$l(\boldsymbol{\beta}, \boldsymbol{\alpha}; \hat{\mathbf{b}}, \mathbf{y}) = \sum_{i=1}^N \left[ \sum_{j=1}^{n_i} \log \left[ p \left( y_{ij} | \hat{\mathbf{b}}_i \right) \right] + n_i \log \left[ p \left( \hat{\mathbf{b}}_i \right) \right] \right. \\ \left. + \frac{1}{2} n_{b_i} \log(2\pi) - \log \left| -\frac{1}{2} n_i f'' \left( \boldsymbol{\beta}, \boldsymbol{\alpha}; \hat{\mathbf{b}}_i \right) \right| \right]$$

where  $n_{b_i}$  is the dimension of the random effect  $\mathbf{b}_i$

- ▶ Maximizing the above log-likelihood function yields the MLE of the Laplace method parameters.
- ▶ This is the default in glmer. The function finds the MLEs of  $\boldsymbol{\beta}$  and  $\boldsymbol{\alpha}$ , using a conditional maximum of  $\mathbf{b}$  (obtained as explained in the following slides).

# Penalized Iterated Reweighted Least Squares (PIRLS)

- For the Laplace approximation, we need to compute the conditional modes of the random effects:

$$\tilde{\mathbf{b}}(\boldsymbol{\beta}, \boldsymbol{\alpha}) = \operatorname{argmax}_b \{p(y|\boldsymbol{\beta}, \mathbf{b}) f(\mathbf{b}, \mathbf{D}(\boldsymbol{\alpha}))\}$$

Conditional modes are calculated using Penalized Iterated Reweighted Least Squares (PIRLS).

- This will be very similar to the IRLS procedure used in the GLM case except that we will add a penalty term:

$$\tilde{\mathbf{b}}(\boldsymbol{\beta}, \boldsymbol{\alpha}) = \operatorname{argmax}_b \left\{ \log p(y|\boldsymbol{\beta}, \mathbf{b}) - \frac{\mathbf{b}^T \mathbf{D}^{-1}(\boldsymbol{\alpha}) \mathbf{b}}{2} \right\}$$

- To de-clutter notation write  $\mathbf{b}$  and  $\boldsymbol{\beta}$  in a single vector  $\mathbf{B}^T = (\mathbf{b}^T, \boldsymbol{\beta}^T)$ , and define corresponding model matrix and precision matrix

$$\mathcal{X} = (\mathbf{Z}, \mathbf{X}) \text{ and } \mathbf{S} = \begin{bmatrix} \mathbf{D}(\boldsymbol{\alpha})^{-1} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix}$$



- ▶ Then, one can show that the Hessian of the penalized likelihood can be written as  $-\mathcal{X}^\top \mathbf{W} \mathcal{X} / \phi - \mathbf{S}$ , where  $\mathbf{W}$  is the diagonal weight matrix (Fisher or full Newton) implied by  $\hat{\boldsymbol{\mu}}$  (e.g. Fisher weights in GLM:  $\mathbf{W} = \text{diag}(w_i)$  where  $w_i = 1 / \left\{ g'(\mu_i)^2 V(\mu_i) \right\}$ )
- ▶ Let  $\mathbf{G} = \text{diag} \{ g'(\mu_i) / \alpha(\mu_i) \}$ . Then, a single Newton update step has the form

$$\mathcal{B}^{[k+1]} = \mathcal{B}^{[k]} + (\mathcal{X}^\top \mathbf{W} \mathcal{X} + \phi \mathbf{S})^{-1} \left\{ \mathcal{X}^\top \mathbf{W} \mathbf{G} (\mathbf{y} - \hat{\boldsymbol{\mu}}) - \phi \mathbf{S} \mathcal{B}^{[k]} \right\}$$

- ▶ Substituting  $\mathcal{B}^{[k]} = (\mathcal{X}^\top \mathbf{W} \mathcal{X} + \phi \mathbf{S})^{-1} (\mathcal{X}^\top \mathbf{W} \mathcal{X} + \phi \mathbf{S}) \mathcal{B}^{[k]}$  and re-arranging, then we obtain

$$\mathcal{B}^{[k+1]} = (\mathcal{X}^\top \mathbf{W} \mathcal{X} + \phi \mathbf{S})^{-1} \mathcal{X}^\top \mathbf{W} \left\{ \mathbf{G} (\mathbf{y} - \hat{\boldsymbol{\mu}}) + \mathcal{X} \mathcal{B}^{[k]} \right\}$$

# Penalized Iterated Reweighted Least Squares (PIRLS)

- The last expression can be seen as the minimizer of the penalized weighted least squares objective

$$\|\mathbf{z} - \mathcal{XB}\|_W^2 + \phi \mathcal{B}^\top \mathbf{SB} = \|\mathbf{z} - \mathbf{X}\boldsymbol{\beta} - \mathbf{Z}\mathbf{b}\|_W^2 + \phi \mathbf{b}^\top \boldsymbol{\psi}_\theta^{-1} \mathbf{b}$$

where  $z_i = g'(\hat{\mu}_i)(y_i - \hat{\mu}_i) + \hat{\eta}_i$  with  $\hat{\eta}_i = g(\hat{\mu}_i)$

# Penalized Iterated Reweighted Least Squares (PIRLS)

- The last expression can be seen as the minimizer of the penalized weighted least squares objective

$$\|\mathbf{z} - \mathcal{X}\mathcal{B}\|_W^2 + \phi\mathcal{B}^\top \mathbf{S}\mathbf{B} = \|\mathbf{z} - \mathbf{X}\boldsymbol{\beta} - \mathbf{Z}\mathbf{b}\|_W^2 + \phi\mathbf{b}^\top \boldsymbol{\psi}_\theta^{-1}\mathbf{b}$$

where  $z_i = g'(\hat{\mu}_i)(y_i - \hat{\mu}_i) + \hat{\eta}_i$  with  $\hat{\eta}_i = g(\hat{\mu}_i)$

- PIRLS algorithm:

1. Initialize  $\hat{\mu}_i = \hat{\eta}_i$
2. Compute pseudo-data  $z_i = g'(\hat{\mu}_i)(y_i - \hat{\mu}_i)/\alpha(\hat{\mu}_i) + \hat{\eta}_i$ , and iterative weights  $w_i = \alpha(\hat{\mu}_i) / \left\{ g'(\hat{\mu}_i)^2 V(\hat{\mu}_i) \right\}$  where  $\alpha(\mu_i)$  is either 1 (Fisher scoring) or  $= 1 + (y_i - \mu_i) \{ V'(\mu_i) / V(\mu_i) + g''(\mu_i) / g'(\mu_i) \}$  (Full Newton-Rahpson)
3. Find  $\hat{\boldsymbol{\beta}}, \hat{\mathbf{b}}$ , to minimize the weighted least squares objective

$$\|\mathbf{z} - \mathbf{X}\boldsymbol{\beta} - \mathbf{Z}\mathbf{b}\|_W^2 + \phi\mathbf{b}^\top \mathbf{D}(\alpha)^{-1}\mathbf{b}$$

4. update  $\hat{\boldsymbol{\eta}} = \mathbf{X}\hat{\boldsymbol{\beta}} + \mathbf{Z}\hat{\mathbf{b}}$  and  $\hat{\mu}_i = g^{-1}(\hat{\eta}_i)$  until convergence

# Penalized quasi-likelihood methods

- ▶ Using Laplace + PIRLS one obtains estimates of  $\beta$ ,  $\alpha$  and  $b$ . However, the nested optimization required to optimize the Laplace approximate profile is somewhat involved and can be computationally costly.
- ▶ As an alternative one can use the PQL to estimate the parameters  $\beta$ ,  $b$  and  $\alpha$  at once using working mixed model. The PQL method estimates the fixed and random effects parameters separately, and then iterate until a convergence criterion is met.

# Penalized quasi-likelihood methods

- First note that, given a current estimate  $\tilde{\mathbf{b}}$  of  $\mathbf{b}$  (e.g., empirical Bayes estimate), one can consider

$$E(\mathbf{y}|\mathbf{b}) = \boldsymbol{\mu} = \mathbf{g}^{-1}(\mathbf{X}'\boldsymbol{\beta} + \mathbf{Z}\mathbf{b}) = \mathbf{g}^{-1}(\boldsymbol{\eta})$$

- Let  $\tilde{\boldsymbol{\beta}}$  and  $\tilde{\mathbf{b}}$  be the current estimates of  $\boldsymbol{\beta}$  and  $\mathbf{b}$ . According to Wolfinger and O'Connell (1993), the first-order Taylor series expansion of  $\boldsymbol{\mu}$  about  $\tilde{\boldsymbol{\beta}}$  and  $\tilde{\mathbf{b}}$  yields approximation of mean

$$\mathbf{g}^{-1}(\boldsymbol{\eta}) = \mathbf{g}^{-1}(\tilde{\boldsymbol{\eta}}) + \tilde{\boldsymbol{\Delta}}\mathbf{X}'(\boldsymbol{\beta} - \tilde{\boldsymbol{\beta}}) + \tilde{\boldsymbol{\Delta}}\mathbf{Z}'(\mathbf{b} - \tilde{\mathbf{b}}) + \textit{error}$$

where  $\tilde{\boldsymbol{\Delta}}$  is the diagonal matrix of the derivatives of the conditional mean evaluated at the expansion, mathematically defined by  $\tilde{\boldsymbol{\Delta}} = \left[ \frac{\partial \mathbf{g}^{-1}(\boldsymbol{\eta})}{\partial \boldsymbol{\eta}} \right]_{\tilde{\boldsymbol{\beta}}, \tilde{\mathbf{b}}}$

# Penalized quasi-likelihood methods

- By rearranging terms, the following expression is derived:

$$\tilde{\Delta}^{-1} [\mu - g^{-1}(\tilde{\eta})] + X' \tilde{\beta} + Z \tilde{b} = X' \beta + Z b$$

- We can see the left-hand side as the expected value of the pseudo-response variable conditionally on  $b$

$$\tilde{\Delta}^{-1} [y - g^{-1}(\tilde{\eta})] + X' \tilde{\beta} + Z' \tilde{b} \equiv \tilde{y}$$

- $\tilde{y}$  is assumed to be normally distributed, i.e.

$$\tilde{y} = X' \beta + Z' b + \tilde{\varepsilon}$$

with pseudo-error term  $\tilde{\varepsilon} = \tilde{\Delta}^{-1} [\tilde{y} - g^{-1}(\tilde{\eta})]$  mean-zero normal with variance  $\text{var}(\tilde{\varepsilon} | \beta, b) = \phi \tilde{\Delta}^{-1} V_{y_i | \beta, b_i}(\mu_i) \tilde{\Delta}^{-1}$

# Penalized quasi-likelihood methods

- ▶ Then, the estimation of parameters can be performed by following the standard procedures for fitting a linear mixed model: given starting values of  $\beta$ ,  $D$  and  $\phi$  in the marginal likelihood, the empirical Bayes estimates can be computed for  $b$ , thereby generating the pseudo-data  $\tilde{y}$ . Consequently, the approximate linear mixed model can be fitted, yielding updated estimates for  $\beta$ ,  $D$ , and  $\phi$ . The iterative process continues until a convergence criterion is reached.
- ▶ In many cases PQL works surprisingly well, is highly effective and produces estimates very close to the full Laplace approximation based estimates. However, for some types of data it is problematic.
  - ▶ for modelling over-dispersed low-mean count data
  - ▶ for binary data (in standard implementations of the LMMs, the algorithm estimates  $\phi$  even if it is known; for binary data, there is no real information to estimate  $\phi$ )
  - ▶ if there are only a few observations per random effect, the central limit theorem justification for the method is not valid
  - ▶ A final general disadvantage is that we do not have access to the profile likelihood for the model itself, so that LRT and AIC based on those is not possible.