

# A Likelihood Approach to Fitting Nonlinear Mixed-Effects Models to Pharmacokinetic and Pharmacodynamic Data

Douglas Bates

University of Wisconsin - Madison

[<Bates@Wisc.edu>](mailto:Bates@Wisc.edu)

Slides for this presentation are available at  
[lme4.R-forge.R-project.org/slides/](http://lme4.R-forge.R-project.org/slides/)

Midwest Biopharmaceutical Statistics Workshop  
Muncie, Indiana  
May 26, 2010

# Outline

- 1 Introduction
- 2 Statistical theory, applications and approximations
- 3 Model definition
- 4 The penalized least squares problem
- 5 Comparing estimation methods

# Outline

- 1 Introduction
- 2 Statistical theory, applications and approximations
- 3 Model definition
- 4 The penalized least squares problem
- 5 Comparing estimation methods

# Outline

- 1 Introduction
- 2 Statistical theory, applications and approximations
- 3 Model definition
- 4 The penalized least squares problem
- 5 Comparing estimation methods

# Outline

- 1 Introduction
- 2 Statistical theory, applications and approximations
- 3 Model definition
- 4 The penalized least squares problem
- 5 Comparing estimation methods

# Outline

- 1 Introduction
- 2 Statistical theory, applications and approximations
- 3 Model definition
- 4 The penalized least squares problem
- 5 Comparing estimation methods

# Outline

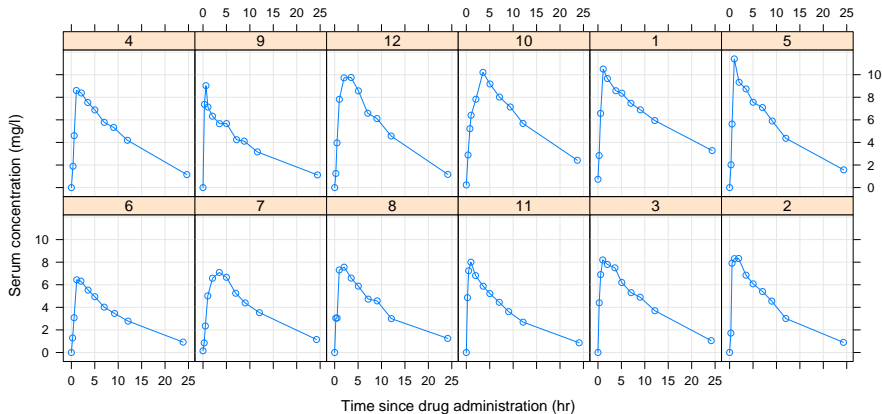
- 1 Introduction
- 2 Statistical theory, applications and approximations
- 3 Model definition
- 4 The penalized least squares problem
- 5 Comparing estimation methods

# Introduction

- Population pharmacokinetic data are often modeled using *nonlinear mixed-effects models* (NLMMs).
- These are *nonlinear* because pharmacokinetic parameters - rate constants, clearance rates, etc. - occur nonlinearly in the model function.
- In statistical terms these are *mixed-effects models* because they involve both *fixed-effects parameters*, applying to the entire population or well-defined subsets of the population, and *random effects* associated with particular experimental or observational units under study.
- Many algorithms for obtaining parameter estimates, usually “something like” the *maximum likelihood estimates* (MLEs), for such models have been proposed and implemented.
- Comparing different algorithms is not easy. Even understanding the definition of the model and the proposed algorithm is not easy.



# An example: Theophylline pharmacokinetics



- These are serum concentration profiles for 12 volunteers after injection of an oral dose of Theophylline, as described in Pinheiro and Bates (2000).

# Modeling pharmacokinetic data with a nonlinear model

- These are longitudinal repeated measures data.
- For such data the time pattern of an individual's response is determined by pharmacokinetic parameters (e.g. rate constants) that occur nonlinearly in the expression for the expected response.
- The form of the nonlinear model is determined by the pharmacokinetic theory, not derived from the data.

$$d \cdot k_e \cdot k_a \cdot C \frac{e^{-k_e t} - e^{-k_a t}}{k_a - k_e}$$

- These pharmacokinetic parameters vary over the population. We wish to characterize typical values in the population and the extent of the variation.
- Thus, we associate random effects with the parameters,  $k_a$ ,  $k_e$  and  $C$  in the nonlinear model.

# Outline

- 1 Introduction
- 2 Statistical theory, applications and approximations
- 3 Model definition
- 4 The penalized least squares problem
- 5 Comparing estimation methods

## Statistical theory and applications - why we need both

- For 30 years, I have had the pleasure of being part of the U. of Wisconsin-Madison Statistics Dept. This year we celebrate the 50th anniversary of the founding of our department by George Box (who turned 90 earlier this year).
- George's approach, emphasizing **both** the theory and the applications of statistics, has now become second-nature to me.
- We are familiar with the dangers of practicing theory without knowledge of applications. As George famously said, "All models are wrong; some models are useful." How can you expect to decide if a model is useful unless you use it?
- We should equally be wary of the application of statistical techniques for which we know the "how" but not the "why". Despite the impression we sometimes give in courses, applied statistics is not just a "black box" collection of formulas into which you pour your data, hoping to get back a p-value that is less than 5%. (In the past many people felt that "applied statistics is the use of SAS" but now we know better.)

# The evolving role of approximation

- When Don Watts and I wrote a book on nonlinear regression we included a quote from Bertrand Russell, “Paradoxically, all exact science is dominated by the idea of approximation”. In translating statistical theory to applied techniques (computing algorithms) we almost always use some approximations.
- Sometimes the theory is deceptively simple (maximum likelihood estimates are the values of the parameters that maximize the likelihood, given the data) but the devil is in the details (so exactly how do I maximize this likelihood?).
- Decades of work by many talented people have provided us with a rich assortment of computational approximations and other tricks to help us get to the desired answer - or at least close to the desired answer.
- It is important to realize that approximations, like all aspects of computing, have a very short shelf life. Books on theory can be useful for decades; books on computing may be outmoded in a few years.

## Failure to revisit assumptions leads to absurdities

- Forty years ago, when I took an intro engineering stats class, we used slide rules or pencil and paper for calculations. Our text took this into account, providing short-cut computational formulas and “rules of thumb” for the use of approximations, plus dozens of pages of tables of probabilities and quantiles.
- Today’s computing resources are unimaginably more sophisticated yet the table of contents of most introductory text hasn’t changed.
- The curriculum still includes using tables to evaluate probabilities, calculating coefficient estimates of a simple linear regression by hand, creating histograms (by hand, probably) to assess a density, approximating a binomial by a Poisson or by a Gaussian for cases not available in the tables, etc.
- Then we make up PDF slides of this content and put the file on a web site for the students to download and follow on their laptops during the lecture. Apparently using the computer to evaluate the probabilities or to fit a model would be cheating - you are supposed to do this by hand.

## And what about nonlinear mixed-effects models?

- Defining the statistical model is subtle and all methods proposed for determining parameter estimates use approximations.
- Often the many forms of approximations are presented as different “types” of estimates from which one can pick and choose.
- In 2007-2008 a consortium of pharma companies, the NLMEc, discussed “next generation” simulation and estimation software for population PK/PD modeling. They issued a set of user requirements for such software including, in section 4.4 on estimation

*The system will support but not be limited to the following estimation methods: FO, FOI, FOCE, FOCEI, Laplacian, Lindstrom and Bates, MCMC, MCPM, SAEM, Gaussian quadrature, and nonparametric methods.*

- Note the emphasis on estimation methods (i.e. algorithms). All of these techniques are supposed to approximate the mle's but that is never mentioned.

# Outline

- 1 Introduction
- 2 Statistical theory, applications and approximations
- 3 Model definition**
- 4 The penalized least squares problem
- 5 Comparing estimation methods



## Linear and nonlinear mixed-effects models

- Both linear and nonlinear mixed-effects models, are based on the  $n$ -dimensional response random variable,  $\mathcal{Y}$ , whose value,  $\mathbf{y}$ , is observed, and the  $q$ -dimensional, unobserved random effects variable,  $\mathcal{B}$ .
- In the models we will consider  $\mathcal{B} \sim \mathcal{N}(\mathbf{0}, \Sigma_\theta)$ . The variance-covariance matrix  $\Sigma_\theta$  can be huge but it is completely determined by a small number of *variance-component parameters*,  $\theta$ .
- The conditional distribution of the response,  $\mathcal{Y}$ , is

$$(\mathcal{Y}|\mathcal{B} = \mathbf{b}) \sim \mathcal{N}(\mu_{\mathcal{Y}|\mathcal{B}}, \sigma^2 \mathbf{I}_n)$$

- The conditional mean,  $\mu_{\mathcal{Y}|\mathcal{B}}$ , depends on  $\mathbf{b}$  and on the fixed-effects parameters,  $\beta$ , through a *linear predictor* expression,  $\mathbf{Z}\mathbf{b} + \mathbf{X}\beta$ .
- For a linear mixed model (LMM),  $\mu_{\mathcal{Y}|\mathcal{B}}$  is exactly the linear predictor. For an NLMM the linear predictor determines the parameter values in the nonlinear model function which then determines the mean.

## Transforming to orthogonal random effects

- We never really form  $\Sigma_\theta$ ; we always work with the *relative covariance factor*,  $\Lambda_\theta$ , defined so that

$$\Sigma_\theta = \sigma^2 \Lambda_\theta \Lambda_\theta^\top.$$

Note that we must allow for  $\Lambda_\theta$  to be less than full rank.

- We define a  $q$ -dimensional “spherical” or “unit” random-effects vector,  $\mathbf{U}$ , such that

$$\mathbf{U} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_q), \quad \mathbf{B} = \Lambda_\theta \mathbf{U} \Rightarrow \text{Var}(\mathbf{B}) = \sigma^2 \Lambda_\theta \Lambda_\theta^\top = \Sigma_\theta.$$

- Setting  $\mathbf{U}_\theta = \mathbf{Z} \Lambda_\theta$ , the linear predictor expression becomes

$$\mathbf{Z}\mathbf{b} + \mathbf{X}\boldsymbol{\beta} = \mathbf{Z}\Lambda_\theta \mathbf{u} + \mathbf{X}\boldsymbol{\beta} = \mathbf{U}_\theta \mathbf{u} + \mathbf{X}\boldsymbol{\beta}.$$

where  $\mathbf{U}_\theta$ , like  $\mathbf{Z}_\theta$  is a large, sparse matrix.

## The conditional mode, $\tilde{\mathbf{u}}_{\theta,\beta}$

- Although the probability model is defined from  $(\mathcal{Y}|\mathbf{u} = \mathbf{u})$ , we observe  $\mathbf{y}$ , not  $\mathbf{u}$  (or  $\mathbf{b}$ ) so we want to work with the other conditional distribution,  $(\mathbf{u}|\mathcal{Y} = \mathbf{y})$ .
- The joint distribution of  $\mathcal{Y}$  and  $\mathbf{u}$  is Gaussian with density

$$\begin{aligned} f_{\mathcal{Y},\mathbf{u}}(\mathbf{y}, \mathbf{u}) &= f_{\mathcal{Y}|\mathbf{u}}(\mathbf{y}|\mathbf{u}) f_{\mathbf{u}}(\mathbf{u}) \\ &= \frac{\exp(-\frac{1}{2\sigma^2}\|\mathbf{y} - \boldsymbol{\mu}_{\mathcal{Y}|\mathbf{u}}\|^2)}{(2\pi\sigma^2)^{n/2}} \frac{\exp(-\frac{1}{2\sigma^2}\|\mathbf{u}\|^2)}{(2\pi\sigma^2)^{q/2}} \\ &= \frac{\exp(-[\|\mathbf{y} - \boldsymbol{\mu}_{\mathcal{Y}|\mathbf{u}}\|^2 + \|\mathbf{u}\|^2] / (2\sigma^2))}{(2\pi\sigma^2)^{(n+q)/2}} \end{aligned}$$

- The mode,  $\tilde{\mathbf{u}}_{\theta,\beta}$ , of the conditional distribution  $(\mathbf{u}|\mathcal{Y} = \mathbf{y})$  (also the conditional mean in the case of an LMM) is

$$\tilde{\mathbf{u}}_{\theta,\beta} = \arg \min_{\mathbf{u}} \left[ \|\mathbf{y} - \boldsymbol{\mu}_{\mathcal{Y}|\mathbf{u}}\|^2 + \|\mathbf{u}\|^2 \right]$$

# Outline

- 1 Introduction
- 2 Statistical theory, applications and approximations
- 3 Model definition
- 4 The penalized least squares problem**
- 5 Comparing estimation methods

## Minimizing a penalized sum of squared residuals

- An expression like  $\|\mathbf{y} - \mu_{\mathbf{y}|\mathbf{u}}\|^2 + \|\mathbf{u}\|^2$  is called a *penalized sum of squared residuals* because  $\|\mathbf{y} - \mu_{\mathbf{y}|\mathbf{u}}\|^2$  is a sum of squared residuals and  $\|\mathbf{u}\|^2$  is a penalty on the size of the vector  $\mathbf{u}$ .
- Determining  $\tilde{\mathbf{u}}_{\theta,\beta}$  as the minimizer of this expression is a *penalized least squares* (PLS) problem. For an LMM it is a *penalized linear least squares problem* that can be solved directly (i.e. without iterating). For an NLMM it is a *penalized nonlinear least squares problem*.
- One way to determine the solution in an LMM is to rephrase it as a linear least squares problem for an extended residual vector

$$\tilde{\mathbf{u}}_{\theta,\beta} = \arg \min_{\mathbf{u}} \left\| \begin{bmatrix} \mathbf{y} - \mathbf{X}\beta \\ \mathbf{0} \end{bmatrix} - \begin{bmatrix} \mathbf{U}_\theta \\ \mathbf{I}_q \end{bmatrix} \mathbf{u} \right\|^2$$

This is sometimes called a *pseudo-data* approach because we create the effect of the penalty term,  $\|\mathbf{u}\|^2$ , by adding “pseudo-observations” to  $\mathbf{y}$  and to the predictor.

## The profiled deviance for LMMs

- We can see that  $\tilde{\mathbf{u}}_{\theta, \beta}$  satisfies  $(\mathbf{U}_{\theta}^{\top} \mathbf{U}_{\theta} + \mathbf{I}_q) \tilde{\mathbf{u}}_{\theta, \beta} = \mathbf{U}_{\theta}^{\top} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})$  which we solve using the sparse Cholesky decomposition

$$\mathbf{L}_{\theta} \mathbf{L}_{\theta}^{\top} = \mathbf{P} (\mathbf{U}_{\theta}^{\top} \mathbf{U}_{\theta} + \mathbf{I}_q) \mathbf{P}^{\top}$$

$\mathbf{P}$  is a permutation matrix that has practical importance but does not affect the theory. The matrix  $\mathbf{L}_{\theta}$  is the sparse, lower-triangular factor.

- Let  $r^2(\boldsymbol{\theta}, \boldsymbol{\beta})$  be the minimum penalized residual sum of squares, then  $\ell(\boldsymbol{\theta}, \boldsymbol{\beta}, \sigma | \mathbf{y}) = \log L(\boldsymbol{\theta}, \boldsymbol{\beta}, \sigma | \mathbf{y})$  can be written

$$-2\ell(\boldsymbol{\theta}, \boldsymbol{\beta}, \sigma | \mathbf{y}) = n \log(2\pi\sigma^2) + \frac{r^2(\boldsymbol{\theta}, \boldsymbol{\beta})}{\sigma^2} + \log(|\mathbf{L}_{\theta}|^2)$$

- The conditional estimate of  $\sigma^2$  is

$$\widehat{\sigma^2}(\boldsymbol{\theta}, \boldsymbol{\beta}) = \frac{r^2(\boldsymbol{\theta}, \boldsymbol{\beta})}{n}$$

producing the *profiled deviance*

$$-2\tilde{\ell}(\boldsymbol{\theta}, \boldsymbol{\beta} | \mathbf{y}) = \log(|\mathbf{L}_{\theta}|^2) + n \left[ 1 + \log \left( \frac{2\pi r^2(\boldsymbol{\theta}, \boldsymbol{\beta})}{n} \right) \right]$$

## Profiling the deviance with respect to $\beta$ for LMMs

- In a LMM the deviance depends on  $\beta$  only through  $r^2(\theta, \beta)$  we can obtain the conditional estimate,  $\hat{\beta}_\theta$ , by extending the PLS problem to

$$r^2(\theta) = \min_{\mathbf{u}, \beta} \left[ \|\mathbf{y} - \mathbf{X}\beta - \mathbf{U}_\theta \mathbf{u}\|^2 + \|\mathbf{u}\|^2 \right]$$

with the solution satisfying the equations

$$\begin{bmatrix} \mathbf{U}_\theta^\top \mathbf{U}_\theta + \mathbf{I}_q & \mathbf{U}_\theta^\top \mathbf{X} \\ \mathbf{X}^\top \mathbf{U}_\theta & \mathbf{X}^\top \mathbf{X} \end{bmatrix} \begin{bmatrix} \tilde{\mathbf{u}}_\theta \\ \hat{\beta}_\theta \end{bmatrix} = \begin{bmatrix} \mathbf{U}_\theta^\top \mathbf{y} \\ \mathbf{X}^\top \mathbf{y} \end{bmatrix}$$

- The profiled deviance, which is a function of  $\theta$  only, is

$$-2\tilde{\ell}(\theta) = \log(|\mathbf{L}_\theta|^2) + n \left[ 1 + \log \left( \frac{2\pi r^2(\theta)}{n} \right) \right]$$

# Conditional mode and profiled Laplace approximation for NLMMs

- As previously stated, determining the conditional mode

$$\tilde{\mathbf{u}}_{\theta, \beta} = \arg \min_{\mathbf{u}} \left[ \|\mathbf{y} - \boldsymbol{\mu}_{\mathbf{y}|\mathbf{u}}\|^2 + \|\mathbf{u}\|^2 \right]$$

in an NLMM is a penalized nonlinear least squares (PNLS) problem.

- It is a nonlinear optimization problem but a comparatively simple one. The penalty term *regularizes* the optimization.
- The Laplace approximation to the profiled deviance (profiled over  $\sigma^2$ ) is, as before,

$$-2\tilde{\ell}(\boldsymbol{\theta}, \boldsymbol{\beta}|\mathbf{y}) = \log(|\mathbf{L}_{\boldsymbol{\theta}}|^2) + n \left[ 1 + \log \left( \frac{2\pi r^2(\boldsymbol{\theta}, \boldsymbol{\beta})}{n} \right) \right]$$

where  $\mathbf{L}_{\boldsymbol{\theta}}$  is the sparse Cholesky factor evaluated at the conditional mode.

- The motivation for this approximation is that it replaces the conditional distribution,  $(\mathbf{U}|\mathbf{Y} = \mathbf{y})$ , for parameters  $\boldsymbol{\beta}$ ,  $\boldsymbol{\theta}$  and  $\sigma$ , by a multivariate Gaussian approximation, *evaluated at the mode*.



# Laplace approximation and adaptive Gauss-Hermite quadrature

- The Laplace approximation

$$-2\tilde{\ell}(\boldsymbol{\theta}, \boldsymbol{\beta}|\mathbf{y}) = \log(|\mathbf{L}_{\boldsymbol{\theta}}|^2) + n \left[ 1 + \log \left( \frac{2\pi r^2(\boldsymbol{\theta}, \boldsymbol{\beta})}{n} \right) \right]$$

is a type of *smoothing objective* consisting of two terms:  $n \left[ 1 + \log \left( \frac{2\pi r^2(\boldsymbol{\theta}, \boldsymbol{\beta})}{n} \right) \right]$ , which measures fidelity to the data, and  $\log(|\mathbf{L}_{\boldsymbol{\theta}}|^2)$ , which measures the complexity of the model.

- For models with a simple structure for the random effects (the matrices  $\boldsymbol{\Sigma}_{\boldsymbol{\theta}}$  and  $\boldsymbol{\Lambda}_{\boldsymbol{\theta}}$  are block diagonal consisting of a large number of small blocks) a further enhancement is to use *adaptive Gauss-Hermite quadrature*, requiring values of the RSS at several points near  $\tilde{\boldsymbol{u}}_{\boldsymbol{\theta}, \boldsymbol{\beta}}$
- Note that the modifier *adaptive*, meaning evaluating at the conditional mode, is important. Gauss-Hermite quadrature without first determining the conditional mode is not a good idea.

# Outline

- 1 Introduction
- 2 Statistical theory, applications and approximations
- 3 Model definition
- 4 The penalized least squares problem
- 5 Comparing estimation methods**

## Consequences for comparisons of methods

- We should distinguish between an algorithm, which is a sort of a black box, and a criterion, such as maximizing the likelihood (or, equivalently, minimizing the deviance).
- The criterion is based on the statistical model and exists outside of any particular implementation or computing hardware. It is part of the theory, which has a long shelf life.
- A particular approximation, algorithm and implementation has a short shelf life.
- I claim it does not make sense to regard the FO, FOI, ... methods as producing well-defined types of “estimates” in the same sense that maximum likelihood estimates, or maximum *a posteriori* estimates are defined.
- If you use a criterion to define an estimation method then implementations should be compared on the basis of that criterion, not on something like mean squared error.