

## Simple fitting of subject-specific curves for longitudinal data

M. Durbán<sup>1,\*†</sup>, J. Harezlak<sup>2</sup>, M.P. Wand<sup>3</sup>, R.J. Carroll<sup>4</sup>

<sup>1</sup> *Department of Statistics, Universidad Carlos III de Madrid, Leganés, 28199 Madrid, Spain.*

<sup>2</sup> *Department of Biostatistics, Harvard School of Public Health, 665 Huntington Avenue, Boston, Massachusetts 02115, USA.*

<sup>3</sup> *Department of Statistics, University of New South Wales, Sydney, NSW, 2052, Australia.*

<sup>4</sup> *Department of Statistics, Texas A&M University, College Station, Texas 77842, USA.*

### SUMMARY

We present a simple semiparametric model for fitting subject-specific curves for longitudinal data. Individual curves are modeled as penalized splines with random coefficients. This model has a mixed model representation, and it is easily implemented in standard statistical software. We conduct an analysis of the long-term effect of radiation therapy on the height of children suffering from acute lymphoblastic leukemia using penalized splines in the framework of semiparametric mixed effects models. The analysis revealed significant differences between therapies and showed that the growth rate of girls in the study cannot be fully explained by the group-average curve and that individual curves are necessary to reflect the individual response to treatment. We also show how to implement these models in **S-PLUS** and **R** in the appendix. Copyright © 2004 John Wiley & Sons, Ltd.

**KEY WORDS:** Linear mixed models; Restricted likelihood ratio tests; Penalized splines; Acute lymphoblastic leukemia.

---

\*Correspondence to: Department of Statistics, Universidad Carlos III de Madrid, Leganés, 28199 Madrid, Spain

†E-mail:mdurban@est-econ.uc3m.es

## 1. INTRODUCTION

Longitudinal data arise frequently in many medical and biological applications. They generally involve a collection of data at different time points for several subjects, and they are characterized by the dependence of repeated observations over time within the same subject. The basic random-effects models for longitudinal data represent each individual as the sum of a population mean which depends on time and is modeled as fixed effect, and a low degree polynomial with random coefficients to model the individual variation. This approach yields a mixed effects model which provides a flexible framework to analyze these type of data [1]. In many situations this parametric assumption is not appropriate.

In many situations, the objective of a longitudinal study is to describe how the response variable is affected by time and other covariates, and the features of the individual profiles. The time course is often too complicated to be model parametrically, this is the reason why in recent years there has been increasing interest in nonparametric analysis of longitudinal data and more specifically in nonparametric subject-specific curves. Early work in this context [2, 3] proposed a semiparametric mixed model for longitudinal data and used different types of smoothers (kernels, smoothing splines, etc.) to estimate the mean population curve, but random effects were modeled by parametric functions. Zhang et al. [4] extended this work by considering a more general class of semiparametric stochastic models by accounting for the within-subject correlation using a stationary or non-stationary Gaussian process, though they did not consider smooth curves for individual subjects. Brumback and Rice [5] modeled both population mean and subject-specific curves nonparametrically with smoothing splines

and used their mixed model representation to present a unified approach for these types of models. However, they ran into computational problems because they assumed fixed slopes and intercepts for the subject-specific curves. Rice and Wu [6] partially solved this problem by modeling individual curves as spline functions with random coefficients. However, in their low-rank spline basis approach the number and location of the knots used to construct the basis became an important issue. Consequently the fit of their models involved the use of some selection criteria to choose these parameters. More recently Guo [7] took a functional data analysis approach by introducing functional random effects which are modeled as realizations of a zero-mean stochastic process. He also used the connection between smoothing splines and mixed models for fitting and estimation of his model. However, Guo [7] also faced computational problems due to large matrices (since smoothing splines use as many knots as data points), and consequently developed a sequential estimation procedure using Kalman filtering [8].

Our approach is a trade-off between spline regression (too dependent on number and position of knots) and smoothing splines (too computationally intensive with large data sets). We use low-rank smoothers (as in [6]) with a penalty approach [9]. We also use the equivalence between a penalized smoother and the optimal predictor in a mixed model [10, 11, 12] to present a unified approach for model estimation. The penalty approach relaxes the importance of the number and location of the knots and the use of a low-rank smoother solves the computational problems of other approaches when analyzing large data sets. The mixed model approach allows one to use existing statistical software and simplify the fit of otherwise complicated models.

The paper is organized as follows: Section 2 presents a general semiparametric mixed model for longitudinal data based on penalized splines and their correspondence with the optimal

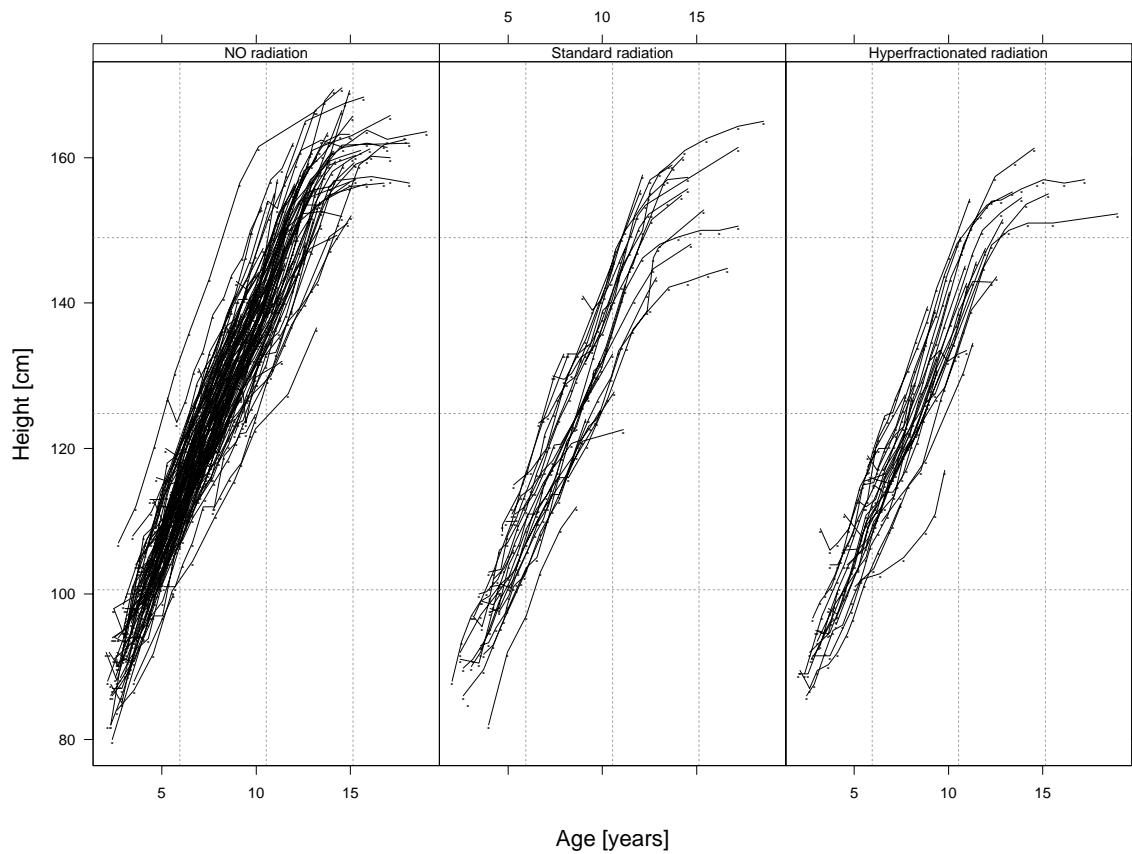


Figure 1. Height of girls over time for each treatment received

predictor in the mixed models. Details on estimation and inference about the parameters are given in Section 3. The models proposed are illustrated in Section 4 with the analysis of data from children suffering from acute lymphoblastic leukemia (ALL). The paper concludes with a brief discussion in Section 5 and an appendix where we show in detail how to implement these models in S-PLUS and R.

## 2. SEMIPARAMETRIC MIXED MODELS FOR LONGITUDINAL DATA

Mixed models are the most widely used method for the analysis of longitudinal data. However, the parametric assumption in the linear mixed model may not always be appropriate, and in many longitudinal studies the response should be modeled as a non-linear function of time for each individual.

We propose the penalized spline model to model the deviation of each subject curve from the population average. The correspondence between the penalized spline smoother and the optimal predictor in a mixed model allows us to take advantage of the methodology and software existent for mixed model analysis, and makes possible a simple implementation of otherwise complicated models. We will illustrate the models proposed with the analysis of the data collected on children suffering from acute lymphoblastic leukemia (ALL). Obesity and short stature are common late effects for childhood ALL and treatments are aimed to minimize the side effects without compromising efficacy. In one of the clinical trials carried out at Dana Farber Cancer Institute (Boston, USA) (see [13] for a detailed description) a total of 618 children were treated between November 1987 and December 1995 with three different central nervous system therapies: intrathecal therapy alone (no radiation), intrathecal therapy with conventional cranial radiation, and intrathecal therapy with twice daily radiation. Measurements on height and weight were taken at diagnosis and approximately every 6 months thereafter. Previous studies on the effects of cranial radiation on height suggested that radiation contributed to decreased expected height, since cranial radiation has been associated with the development of growth hormone deficiency.

The purpose of this analysis is to evaluate the long-term effects of treatment on the children height and on the individual growth trajectories.

### 2.1. Random intercept and slope models

Let  $y_{ij}$  denote the height of girl  $i$ ,  $i = 1, \dots, m$  at age  $x_{ij}$ ,  $j = 1, \dots, n_i$ . A starting model for these data could be the linear mixed model proposed by Laird and Ware [1]:

$$y_{ij} = \beta_0 + \beta_1 x_{ij} + U_i + \varepsilon_{ij} \quad (1)$$

where  $U_i \sim N(0, \sigma_U^2)$  and  $\varepsilon_{ij} \sim N(0, \sigma_\varepsilon^2)$ ,  $\beta_0$  represents the overall mean and  $U_i$  is a random intercept for girl  $i$ , which is treated as a random sample from the population of girls and requires just a single parameter,  $\sigma_U^2$ . We can write model (1) in matrix notation as:

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u} + \boldsymbol{\varepsilon} \quad (2)$$

where

$$\mathbf{Y} = \begin{bmatrix} y_{11} \\ \vdots \\ y_{mn_m} \end{bmatrix}, \quad \mathbf{X} = \begin{bmatrix} \mathbf{X}_1 \\ \vdots \\ \mathbf{X}_m \end{bmatrix}, \quad \mathbf{X}_i = \begin{bmatrix} 1 & t_{i1} \\ \vdots & \vdots \\ 1 & t_{in_i} \end{bmatrix},$$

$$\mathbf{Z} = \begin{bmatrix} \mathbf{1}_1 & 0 & \dots & 0 \\ 0 & \mathbf{1}_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \mathbf{1}_m \end{bmatrix}, \quad \mathbf{1}_i = \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix}_{n_i \times 1} \quad \text{and} \quad \boldsymbol{\beta} = [\beta_0, \beta_1]^T$$

Figure 1 plots the height of girls as a function of age for three different treatments. We can see that the linearity assumption is not reasonable for ALL data, and in general, for height of children in this range of age, and so (1) should be at least extended to

$$y_{ij} = f(x_{ij}) + U_i + \varepsilon_{ij} \quad (3)$$

where  $f$  is a smooth function which reflects the overall increasing trend of height along age.

We estimate  $f$  by a penalized spline. Let  $\kappa_1, \dots, \kappa_K$  be a set of distinct knots in the range of

$x_{ij}$  and  $x_+ = \max(0, x)$ . The number of knots  $K$  is fixed and large enough (in this case  $K=40$ ) to ensure the flexibility of the curve. The knots are chosen as quantiles of  $x$  with probabilities  $1/(K+1), \dots, K/(K+1)$ . We use truncated lines as the basis for regression since their simple mathematical form is very useful when formulating complicated models. More complex basis such as B-splines and radial basis functions (with better numerical properties) could also be used to fit these models and the programs presented in the appendix are easily adapted to use any basis for regression. A penalized linear spline model for (3) is:

$$y_{ij} = \underbrace{\beta_0 + \beta_1 x_{ij} + \sum_{k=1}^K u_k (x_{ij} - \kappa_k)_+}_{f(x_{ij})} + U_i + \varepsilon_{ij} \quad U_i \sim N(0, \sigma_U^2) \quad \varepsilon_{ij} \sim N(0, \sigma_\varepsilon^2). \quad (4)$$

For given values of  $\sigma_U^2$  and  $\sigma_\varepsilon^2$ , the estimates of  $(\boldsymbol{\beta}, \mathbf{u})$  are obtained by minimizing the penalized least squares:

$$\sum_{i=1}^m \sum_{j=1}^{n_i} \{y_{ij} - f(x_{ij}) - U_i\}^2 + \lambda \mathbf{u}^T \mathbf{u} + (\sigma_\varepsilon^2 / \sigma_U^2) \mathbf{U}^T \mathbf{U},$$

where the smoothing parameter  $\lambda$  controls the amount of smoothing of  $f$ , i.e., a value of  $\lambda = 0$  is equivalent to ordinary least squares, while with larger values of  $\lambda$ , the term  $\lambda \mathbf{u}^T \mathbf{u}$  looks after the overparametrization of the regression function by placing a penalty on the smoothness of the  $u_k$ , yielding smoother fitted curves. The penalized spline smoother corresponds to the optimal predictor in a mixed model framework assuming  $u_k \sim N(0, \sigma_u^2)$  and  $\lambda = \sigma_\varepsilon^2 / \sigma_u^2$ . Brumback et al. [10], Currie and Durbán [11] and Wand [12], among others, discuss the mixed model representation of penalized splines. This enables us to write (4) as a semiparametric

mixed model similar to (2) where now

$$\mathbf{Z} = \begin{bmatrix} \mathbf{Z}_1 & \mathbf{1}_1 & 0 & \dots & 0 \\ \mathbf{Z}_2 & 0 & \mathbf{1}_2 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \mathbf{Z}_m & \vdots & \vdots & \dots & \mathbf{1}_m \end{bmatrix}, \quad \mathbf{Z}_i = \begin{bmatrix} (t_{i1} - \kappa_1)_+ & \dots & (t_{i1} - \kappa_K)_+ \\ \vdots & \ddots & \vdots \\ (t_{in_i} - \kappa_1)_+ & \dots & (t_{in_i} - \kappa_K)_+ \end{bmatrix},$$

$$\mathbf{u} = [u_1, \dots, u_K, U_1, \dots, U_m]^T \quad \text{and} \quad \mathbf{G} = \text{Cov}(\mathbf{u}) = \begin{bmatrix} \sigma_u^2 \mathbf{I} & 0 \\ 0 & \sigma_U^2 \mathbf{I} \end{bmatrix}.$$

In this model, the deviation of the  $i$ th girl is modeled through a random intercept. This is quite simplistic and will usually not be realistic since it assumes that the growth curves are parallel. A simple extension is to assume that the subject-specific differences are straight lines:

$$y_{ij} = f(x_{ij}) + a_{i1} + a_{i2}x_{ij} + \varepsilon_{ij} \quad \varepsilon_{ij} \sim N(0, \sigma_\varepsilon^2) \quad (a_{i1}, a_{i2})^T \sim N(0, \Sigma) \quad (5)$$

which in matrix notation becomes

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u} + \boldsymbol{\varepsilon},$$

$$\mathbf{Z} = \begin{bmatrix} \mathbf{Z}_1 & \mathbf{X}_1 & \mathbf{0} & \dots & \mathbf{0} \\ \mathbf{Z}_2 & \mathbf{0} & \mathbf{X}_2 & \dots & \mathbf{0} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \mathbf{Z}_m & \mathbf{0} & \mathbf{0} & \dots & \mathbf{X}_m \end{bmatrix}, \quad \mathbf{u} = [u_1, \dots, u_K, a_{11}, a_{12}, \dots, a_{m1}, a_{m2}]^T,$$

$$\mathbf{G} = \text{Cov}(\mathbf{u}) = \begin{bmatrix} \sigma_u^2 \mathbf{I} & 0 \\ 0 & \text{blockdiagonal } \Sigma \\ & & \underset{1 \leq i \leq m}{\phantom{\text{blockdiagonal } \Sigma}} \end{bmatrix}.$$

Another point worth mentioning is that we allow for complex departures from the common linear component, since  $\Sigma$  is an unstructured  $2 \times 2$  matrix. The alternative is to assume that

the subject-specific intercepts and slopes are independent, and hence that  $\Sigma$  is diagonal, and assumption that we do not recommend. In the first instance, if one assumes that  $\Sigma$  is diagonal in the original parameterization  $a_{i1} + a_{i2}x_{ij}$ , then they will not be independent if the  $x$ 's become centered at their mean, as might reasonably be done for numerical stability. Previous papers [14, 7] propose the use of this covariance structure, but the authors do not use it in their examples and do not show how to implement it in the usual statistical packages; in the appendix we show how to implement it in S-PLUS and R.

## 2.2. Subject-specific curves

The most flexible models are those that allow for the subject-specific differences to be non-parametric functions (see for example [15]). This can be done using penalized splines as follows:

$$\begin{aligned} y_{ij} &= f(x_{ij}) + g_i(x_{ij}) + \varepsilon_{ij}, \quad \varepsilon_{ij} \sim N(0, \sigma_\varepsilon^2), \\ g_i(x_{ij}) &= a_{i1} + a_{i2}x_{ij} + \sum_{k=1}^K v_{ik}(x_{ij} - \kappa_k)_+, \quad (a_{i1}, a_{i2})^T \sim N(0, \Sigma) \quad v_{ik} \sim N(0, \sigma_v^2). \end{aligned} \tag{6}$$

This model is an extension of (5) since in that model, the individual trajectories were linear,  $a_{i1} + a_{i2}x_{ij}$ , and in (6) each subject-specific curve has two components: a linear (similar to (5)) and a non-linear part,  $\sum_{k=1}^K v_{ik}(x_{ij} - \kappa_k)_+$ , which allows more flexibility. Both components are random, different from the approach taken by Brumback and Rice [5]. Though model (6)

is complex, it is easily described in the mixed model framework as

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u} + \boldsymbol{\varepsilon} \quad \text{with}$$

$$\mathbf{Z} = \begin{bmatrix} \mathbf{Z}_1 & \mathbf{X}_1 & \mathbf{0} & \dots & \mathbf{0} & \mathbf{Z}_1 & \mathbf{0} & \dots & \mathbf{0} \\ \mathbf{Z}_2 & \mathbf{0} & \mathbf{X}_2 & \dots & \mathbf{0} & \mathbf{0} & \mathbf{Z}_2 & \dots & \mathbf{0} \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots & \vdots & \ddots & \vdots \\ \mathbf{Z}_m & \mathbf{0} & \mathbf{0} & \dots & \mathbf{X}_m & \mathbf{0} & \mathbf{0} & \dots & \mathbf{Z}_m \end{bmatrix}, \quad (7)$$

$$\mathbf{u} = [u_1, \dots, u_K, a_{11}, a_{12}, \dots, a_{m1}, a_{m2}, v_{11}, \dots, v_{mK}]^T,$$

$$\mathbf{G} = \text{Cov}(\mathbf{u}) = \begin{bmatrix} \sigma_u^2 \mathbf{I} & 0 & 0 \\ 0 & \text{blockdiagonal } \boldsymbol{\Sigma}_{1 \leq i \leq m} & \\ 0 & 0 & \sigma_v^2 \mathbf{I} \end{bmatrix}.$$

### 2.3. Factor by curve interactions

One of the purposes of the study carried out with children suffering from ALL was to compare the long-term effects of the three different therapies, so we might be interested in fitting a separate mean curve for children receiving each therapy. To do that we use an interaction model in which a categorical factor interacts with a continuous predictor, so that model (6) can be extended to

$$y_{ij} = f_{z_i}(x_{ij}) + g_i(x_{ij}) + \varepsilon_{ij}$$

$$y_{ij} = \beta_0 + \beta_1 x_{ij} + \sum_{k=1}^K u_k (x_{ij} - \kappa_k)_+ + \sum_{l=2}^L z_{il} (\gamma_{0l} + \gamma_{1l} x_{ij}) + \sum_{l=2}^L z_{il} \left\{ \sum_{k=1}^K w_k^l (x_{ij} - \kappa_k)_+ \right\}$$

$$+ a_{i1} + a_{i2} x_{ij} + \sum_{k=1}^K v_{ik} (x_{ij} - \kappa_k)_+ + \varepsilon_{ij}$$

$$u_k \sim N(0, \sigma_u^2), \quad w_k^l \sim N(0, \sigma_w^2), \quad v_{ik} \sim N(0, \sigma_v^2), \quad \varepsilon_{ij} \sim N(0, \sigma_\varepsilon^2), \quad (8)$$

where  $z_{il} = 1$  if  $z_i = l$  and 0 otherwise for  $l=2,3$ . For simplicity, we have assumed a common variance parameter for all curves, i.e.,  $\text{Var}(w_k^l) = \sigma_w^2$ ,  $l = 2, \dots, L$ . A common variance means

that all curves have equivalent smoothness, but the random effects are independent from function to function, i.e., the curves are different but with the same amount of smoothing. In order for the fixed effects to be identified we need to put constraints on  $\gamma_{jl}$ , we assume  $\gamma_{01} = \gamma_{11} = 0$  (as in [16]) which means that  $\beta_0 + \beta_1 x_{ij} + \sum_{k=1}^K u_k (x_{ij} - \kappa_k)_+$  is the fitted curve for  $l = 1$  and  $\gamma_{0l} + \gamma_{1l} x_{ij} + \sum_{k=1}^K w_k^l (x_{ij} - \kappa_k)_+$  is the difference of the fitted curves for the therapies 2 (standard radiation) and 3 (hyper-fractionated radiation) and therapy 1 (no radiation). The mixed model representation of this model is similar to (6) but now the part of  $\mathbf{Z}$  corresponding to the overall mean is block-diagonal and each block corresponds to the truncated line basis for children receiving each therapy.

### 3. INFERENCE

The linear mixed model representation of penalized splines is the foundation for fitting the models described in Section 2. A standard estimation criterion for variance components is the Restricted Maximum Likelihood (REML) of Patterson and Thompson [17]. For example, given model (6)

$$\begin{aligned} \ell_R(\sigma_u^2, \sigma_v^2, \sigma_\varepsilon^2) = & -\frac{1}{2} \log |\mathbf{V}| - \frac{1}{2} \log |\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X}| \\ & - \frac{1}{2} \mathbf{y}^T (\mathbf{V}^{-1} - \mathbf{V}^{-1} \mathbf{X} (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1}) \mathbf{y}, \end{aligned} \quad (9)$$

where  $\mathbf{V} = \mathbf{Z}\mathbf{G}\mathbf{Z}^T + \sigma_\varepsilon^2 \mathbf{I}$  and  $\mathbf{G}$  is defined in (7). The vector of parameters  $\boldsymbol{\beta}$  and the random coefficient vector  $\mathbf{u}$  can be determined using *best prediction*:

$$\begin{aligned} \hat{\boldsymbol{\beta}} &= (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1} \mathbf{y} \\ \hat{\mathbf{u}} &= \hat{\mathbf{G}} \mathbf{Z}^T \hat{\mathbf{V}}^{-1} (\mathbf{y} - \mathbf{X} \hat{\boldsymbol{\beta}}). \end{aligned}$$

Testing the adequacy of a parametric model against a nonparametric alternative is not straightforward. For example, in model (3) we might be interested in testing whether the function describing the population mean is a line or there is some degree of nonlinearity. This is equivalent to testing

$$H_0 : \sigma_u^2 = 0 \quad \text{vs.} \quad H_1 : \sigma_u^2 > 0.$$

The main problem we face here is that the parameter of interest is on the boundary of the parameter space,  $[0, \infty)$ , so the restricted likelihood ratio statistic

$$RLRT = \sup_{H_1} REL(\boldsymbol{\beta}, \sigma_\varepsilon^2, \sigma_U^2, \sigma_u^2) - \sup_{H_0} REL(\boldsymbol{\beta}, \sigma_\varepsilon^2, \sigma_U^2, \sigma_u^2) \quad (10)$$

can not be compared with a  $\chi_1^2$ . Self and Liang [18] and Stram and Lee [19] discussed the asymptotic distribution of RLRT and showed that under the assumption that  $\mathbf{y}$  can be partitioned into independent subvectors and the number of subvectors tends to infinity, (10) has a  $\frac{1}{2}\chi_q^2 + \frac{1}{2}\chi_{q+1}^2$  asymptotic distribution, where  $q$  is the number of fixed effects under the null hypothesis. However, this assumption does not hold under the alternative hypothesis in this type of semiparametric mixed models, and the chi-squared mixture approximation can be poor [20]. Crainiceanu and Ruppert [21] derived the case of testing polynomial regression against a general alternative modeled by penalized splines for one variance component and Crainiceanu et al. [22] studied the case where there are several variance components. These authors also suggest the use of simulation to determine the null distribution of the likelihood ratio test statistic. The idea is to estimate the model parameters under the null hypothesis, then simulate the distribution of the likelihood ratio test under the null model at the parameters. Crainiceanu et al. [22] give fast simulation algorithms in some cases, however, the complexity of these algorithms increases linearly with the number of subjects and the fitting of complex models to a large number of simulated datasets can become computationally infeasible. The

complexity of the model proposed in this paper, and the large number of observations in the dataset makes it difficult to implement the methods proposed above. As a guide line, we will compare the RLRT with the chi-squared mixture approximation.

#### 4. APPLICATION TO THE ACUTE LYMPHOBLASTIC LEUKEMIA DATA

We used the semiparametric mixed models described in Section 2 to analyze the Acute Lymphoblastic Leukemia (ALL) data. We concentrate on analysis of 197 girls diagnosed with ALL between 2 and 9 years of age. Height was measured at different times and a total of 1988 observations were obtained. The number of observations per girl ranged from 1 to 21. Two nested models with 5 and 6 variance components were fitted to the data, namely models (5) extended with factor by curve interaction and model (8) :

$$y_{ij} = f_{z_i}(x_{ij}) + a_{i1} + a_{i2}x_{ij} + \varepsilon_{ij},$$

$$y_{ij} = f_{z_i}(x_{ij}) + g_i(x_{ij}) + \varepsilon_{ij},$$

where  $y_{ij}$  is the height in centimeters of the  $i$ -th girl at age  $j$  (in years), for  $i = 1, \dots, 197$  and  $j$  between 1 and 21,  $f_1$  is the group-average curve for girls receiving intrathecal therapy alone,  $f_2$  for girls receiving intrathecal therapy with conventional cranial radiation and  $f_3$  for girls receiving intrathecal therapy with twice daily radiation,  $a_{i1}$  and  $a_{i2}$  are random intercepts and slopes respectively, and  $g_i(x_{ij})$  is the subject-specific deviation of the  $i$ -th girl from the group-average curve. We used REML for variance component estimation and the function `lme()` implemented in S-PLUS 2000 and R 1.8 to fit the model. Details can be found in the appendix.

In our analysis, we are interested in estimating mean treatment effects and the individual responses to treatment. In model (8), we assume that curves describing the effect of treatment

have equivalent smoothness. We can test this assumption which is equivalent to the hypothesis:

$$H_0 : \sigma_{w_1}^2 = \sigma_{w_2}^2 = \sigma_{w_3}^2 = \sigma_w^2.$$

The appropriate way to proceed would be to use parametric bootstrap to obtain the distribution of the likelihood ratio test, however, as we mentioned above, the computational time needed to fit the null model to a large number of simulated data sets (between 10,000 and 100,000) would make the bootstrap method infeasible. The distribution theory or Monte Carlo approaches are still an open problem. As a guideline, we compare the  $-2\log(\text{RLRT})=0.7719$  with 1.642 which is the 90th percentile of  $\frac{1}{2}\chi_0^2 + \frac{1}{2}\chi_1^2$  (the distribution of  $-2\log(\text{RLRT})$  under the assumption of independent  $\mathbf{y}$ 's). This result suggests that we do not need separate variance components for each curve. Figure 2 (left) shows the estimates of the population curves for all three groups in model (8). It can be seen that all groups have similar height patterns, but the girls not receiving radiation (treatment 1) are taller than girls in the other two groups. In particular, the group not receiving radiation seems to be significantly taller when they reach adolescence. To compare the three average-curves we refit the model with one common average curve, the null hypothesis corresponding to

$$H_0 : \gamma_{jl} = 0 \quad j = 0, 1 \quad l = 1, 2, 3 \quad \text{and} \quad \sigma_w^2 = 0$$

for  $\gamma_{jl}$  and  $\sigma_w^2$  defined in (8). We compare the  $-2\log(\text{RLRT})=30.11$  with the 90th percentile of  $\frac{1}{2}\chi_2^2 + \frac{1}{2}\chi_3^2$ , 5.528. This result shows a high degree of statistical significance, implying that the height of girls is affected by the treatment received. This result supports the findings that associate cranial radiation with growth hormone deficiency [24, 25]. Crainiceanu and Ruppert [21] test for fixed and random effects simultaneously in models with one variance component, further research needs to be done to extend this work to the present case where

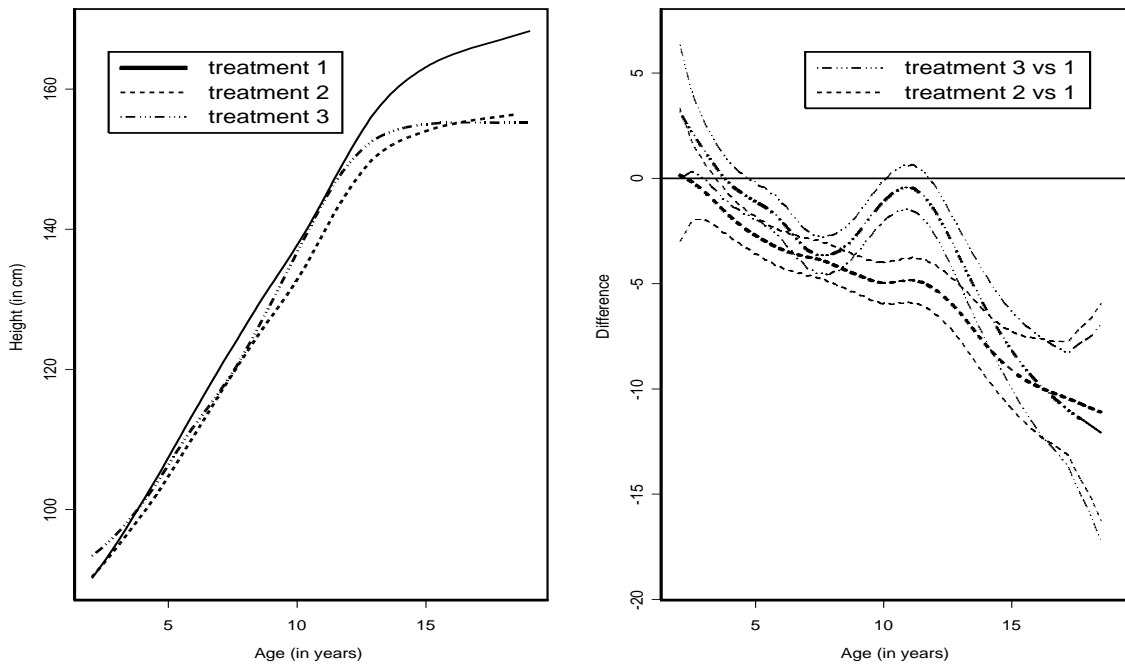


Figure 2. Estimates of the population mean response curves (left) and contrast curves with pointwise confidence bands (right)

there are several variance components. In the right panel of Figure 2, we present the contrast curves,  $\hat{f}_2(x_{ij}) - \hat{f}_1(x_{ij})$  and  $\hat{f}_3(x_{ij}) - \hat{f}_1(x_{ij})$ . The plot shows how the height of girls receiving conventional cranial radiation (treatment 2) is lower, at all ages than the height of girls not receiving radiation (treatment 1). The method presented in Ruppert et al. [15] could be extended to compute simultaneous confidence bands for these curves in an efficient way. The software package ASREML [27] also fits these models and handles confidence bands around curves.

To test whether or not the individual response to treatment is linear we compare models (5)

and (6). Figure 3 shows the estimates of random effects corresponding to the subject-specific curves in model (8). It shows that the between-girls variation is considerable and that a linear random effect to describe the within-subject variation would not be appropriate. Comparing models (5) and (6) (extended with factor by curve interaction) is equivalent to testing:

$$H_0 : \sigma_v^2 = 0 \quad \text{vs.} \quad H_1 : \sigma_v^2 > 0.$$

Again, it would be preferable to use Monte Carlo simulations, but given the complexity of the model we compare the value of  $-2 \log(\text{RLRT})=202.99$  to the 90th percentile of a  $\frac{1}{2}\chi_0^2 + \frac{1}{2}\chi_1^2$ , 1.642. This result indicates that the deviation of the each girl from the population average needs to be modeled nonparametrically. Figure 4 shows individual and treatment group curves for six girls. The average group curves do not follow the response of individual girls and are different from the individual curves, showing the loss of information about individual trajectories when subject-specific curves are not included in the model. The apparent near interpolation of the fitted curves is due to the difference in the scale of the height trend compared with the error variance. One last remark is the fact that not accounting for the individual variation correctly can have an impact on the comparison of the average-curves for each therapy. In this example, there was no significant difference between the average-curves for each group when the model fitted only included a random intercept and did not include subject-specific trajectories.

## 5. CONCLUSIONS

We have provided a flexible and simple method of fitting individual curves in longitudinal studies. The mixed model representation of penalized splines allows one to take advantage of the existing methodology for mixed model analysis and the use of software such as PROC MIXED

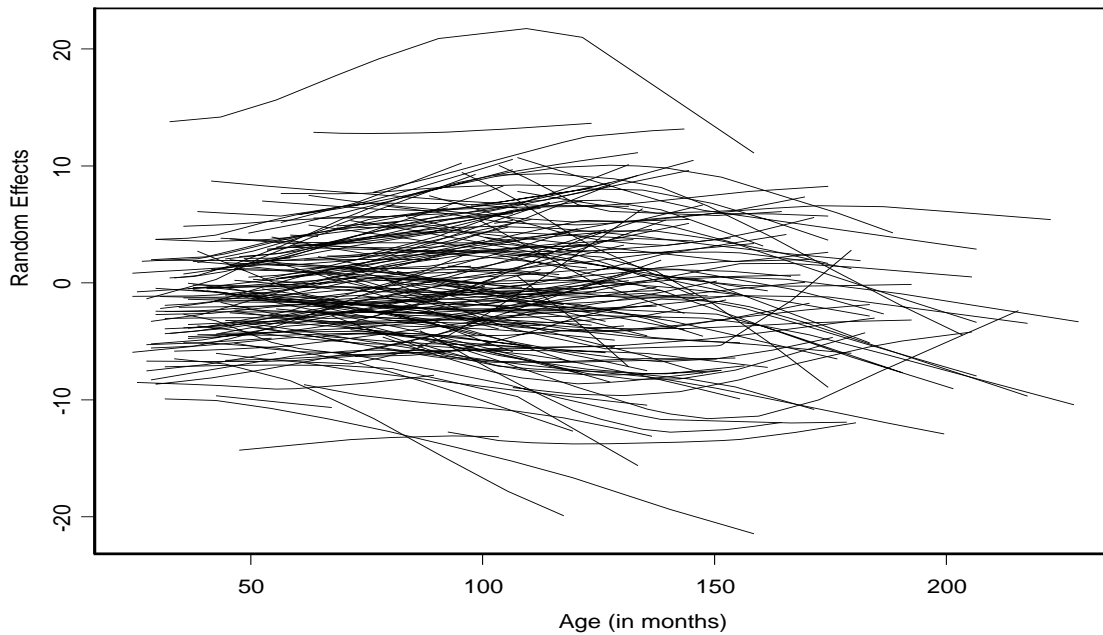


Figure 3. Estimates of random effects for each girl

in SAS and `lme()` in S-PLUS and R. Our use of low-rank smoothers with penalties and random subject-specific curves solves the computational problems of previous approaches [5, 7] based on smoothing splines and reduces significantly the time needed to fit the models. Our approach allows a fast fit of complex models to longitudinal problems with large number of individuals, the fit of model (6) with 197 subjects and 1988 observations took less than a minute on a 1.20 GHz Pentium III PC.

Our analysis of ALL data indicates that the growth of girls who did not receive radiation was not slowed down by the therapy. In our data, girls not receiving radiation were on average taller than girls receiving standard or twice daily radiation, and this difference increased when girls

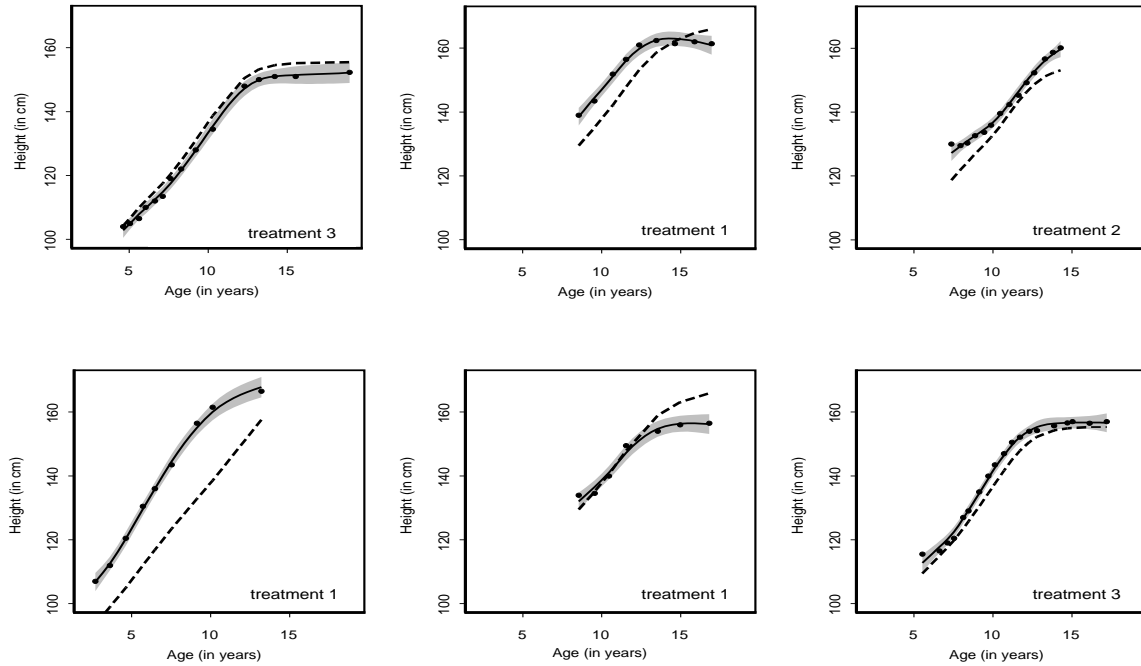


Figure 4. Plot of individual curves (solid lines) with 95% confidence intervals together with treatment average group curves (dashed lines) and observed heights for comparison.

reached puberty. The flexibility of these nonparametric models captured the relative growth effects of the three different therapies on each girl, showing that the growth rate of the girls in the study cannot be fully explained by the group-average curve and that a linear function would not be appropriate to describe the long-term effect of therapy on height. The models presented here can be easily extended to more general models, for example, to account for correlation among errors.

## APPENDIX: IMPLEMENTATION IN STATISTICAL SOFTWARE

Statistical packages such as **S-PLUS**, **R** and **SAS** have generic functions which fit linear mixed-effects models and allow for nested random effects. The **SAS** procedure **PROC MIXED** and the **lme()** function in **S-PLUS** and **R** allow the fitting of complex mixed models and so, obtain easily estimates of fixed and random effects and variance components. Below we describe the fitting of models (3), (6) and (8) to **ALL** data. This dataset is confidential, so it cannot be made available, however, a simulated dataset and the programs to fit these data with **SAS**, **S-PLUS** and **R** are available on request.

## I. Random slope and intercept models

We start by giving some code to set up the basic inputs of the programs. The response variable, factors and independent variables are part of a **S-PLUS** data set, **ALL**:

```
attach(ALL)
y <- ALL$height
time <- ALL$age
treatment <- ALL$xrtrand
subject <- ALL$child
```

We follow Ruppert [26] to set the number and location of the knots used to compute the basis of the penalized spline for the overall mean:

```
K <- max(5,min(floor(length(unique(time)))/4),40)
knots <- quantile(unique(time),seq(0,1,length=K+2))[-c(1,K+2)]
```

Set up design matrices and truncated lines basis:

```
X <- model.matrix(y ~ time)
Z <- outer(time,knots,"-")
Z <- Z*(Z>0)
n <- length(y)
```

The fit of model (3),  $y_{ij} = f(x_{ij}) + U_i + \varepsilon_{ij}$  is given in two line commands:

```
Id <- factor(rep(1,length(y)))
fit <- lme(y~time,random=list(Id=pdIdent(~Z-1),subject=pdIdent(~1)))
```

`Id=` indicates that there is no grouping structure for the design matrix  $Z$ , i.e., there is a single curve for all therapy groups; `pdIdent` specifies that the variance structure of the random effects is a multiple of the identity, and `subject=pdIdent(~1)` indicates that each child has a different random coefficient (a total of 197), but with a single variance component for all of them. The estimated variance components and estimates of fixed and random effects are:

```
sig.sq.hat <- fit$sigma^2
sig.sq.U.hat <- sig.sq.hat*exp(2*unlist(fit$modelStruct)[1])
sig.sq.u.hat <- sig.sq.hat*exp(2*unlist(fit$modelStruct)[2])
beta.hat <- fit$coeff$fixed
u.hat <- unlist(fit$coeff$random)
f.hat <- X%*%beta.hat+Z%*%u.hat[1:ncol(Z)]
d <- dim(fit$fitted)
fitted.val <- fit$fitted[,d]
```

Model (5),  $y_{ij} = f(x_{ij}) + a_{i1} + a_{i2}x_{ij} + \varepsilon_{ij}$  fits a specific slope and intercept for each girl:

```
fit <- lme(y~time,random=list(Id=pdIdent(~Z-1),subject=pdSymm(~time)))
```

The command `subject=pdSymm(~time)` specifies a  $2 \times 2$  symmetric positive-definite matrix covariance structure for the random intercept and slope,  $a_{i1} + a_{i2}x_{ij}$   $(a_{i1}, a_{i2})^T \sim N(0, \Sigma)$ , identical but separate for each subject.

### *I.1. R code*

```
library(nlme)

Z.block<-list(list(Id=pdIdent(~Z-1)),list(subject=pdIdent(~1)))

Z.block<-unlist(Z.block,recursive=FALSE)

data.fr <- groupedData( y ~ X[,-1] | rep(1,length=length(y)),data =
                        data.frame(y,X,Z,subject))

fit <- lme(y~X[,-1],data=data.fr,random=Z.block)
```

For model Model (5),

```
Z.block<-list(list(Id=pdIdent(~Z-1)),list(subject=pdSymm(~time)))

Z.block<-unlist(Z.block,recursive=FALSE)

fit <- lme(y~X[,-1],data=data.fr,random=Z.block)
```

## II. Subject-specific curves

We now set up truncated lines basis for subject-specific curves. The number of observations within subject is at most 21 and the number of individuals is 197, now we will use 10 knots (instead of 40) to construct the basis for each subject.

```
K.subject <- 10
```

```
knots.subject <- quantile(unique(time),seq(0,1,length=K.subject+2)
                          )[-c(1,K.subject+2)]
Z.subject <- outer(time,knots.subject,"-")
Z.subject <- Z.subject*(Z.subject>0)
```

The fit of model (6),  $y_{ij} = f(x_{ij}) + g_i(x_{ij}) + \varepsilon_{ij}$  is:

```
fit <- lme(y~time,random=list(Id=pdIdent(~Z-1),subject=pdSymm(~time),
                             subject=pdIdent(~Z.subject-1)))
```

`subject=pdIdent(~Z.subject-1)` specifies a common diagonal covariance matrix for the deviations from linearity,  $\sum_{k=1}^K v_k(x_{ij} - \kappa_k)_+ \quad v_k \sim N(0, \sigma_v^2)$ .

### II.1. R code

```
Z.block<-list(list(Id=pdIdent(~Z.total-1)),list(case=pdSymm(~time)),
              list(case=pdIdent(~Z.subject-1)))
Z.block<-unlist(Z.block,recursive=FALSE)
data.fr <- groupedData( y ~ X[,-1] |rep(1,length=length(y)),
                       data = data.frame(y,X,Z,Z.subject,case))
fit <- lme(y~X[,-1],data=data.fr,random=Z.block)
```

## III. Factor by curve interactions

We know fit model (8),  $y_{ij} = f_{z_i}(x_{ij}) + g_i(x_{ij}) + \varepsilon_{ij}$ , in order to fit a separate curve mean curve for each therapy group, for simplicity, we assume a common variance parameter for all curves:

```
options(contrasts=c("contr.treatment","contr.poly"))
fit <- lme(y~treatment*time,random=list(treatment=pdIdent(~Z-1),
      subject = pdSymm(~time),subject=pdIdent(~Z.subject-1)))
```

`options(contrasts=c("contr.treatment","contr.poly"))` specifies the constraints used to ensure the identifiability of the model. We choose to set the first level of each factor included in the model equal to 0, `y~treatment*time` fits a separate fixed slope and intercept for the height of children receiving each therapy and `treatment=pdIdent(~Z-1)` specifies a common variance parameter, i.e. a common smoothing parameter, but the shape of the fitted curves is different for each group, `treatment=` indicates that matrix  $Z$  is to be split into blocks, with each block corresponding to a different group (since  $Z$  has 40 columns we have a total of 120 random coefficients for group curves).

### III.1. R code

```
X <- model.matrix(y ~ treatment*time)
Z.block<-list(list(treatment=pdIdent(~Z.total-1)),list(case=
      pdSymm(~time)),list(case=pdIdent(~Z.subject-1)))
Z.block<-unlist(Z.block,recursive=FALSE)
data.fr <- groupedData( y ~ X[,-1] |rep(1,length=length(y)),
      data = data.frame(y,X,Z,Z.subject,case))
fit<- lme(y~X[,-1],data=data.fr,random=Z.block)
```

One of the advantages of using `lme()` function in S-PLUS and R is that we do not need to create the full matrix  $Z$  as described in Section 2. We only need to work with a matrix of size  $1988 \times 40$  or  $1988 \times 10$  instead of a matrix of size  $1988 \times 2484$  (since we have 2484 random effects

in this last model: 120 for factor by group interaction, 394 for random slopes and intercepts and 1970 for departures from linearity), working with so large matrices would be almost impossible. Other packages such as SAS creates full Z matrix, in this case, the number of knots for the individual curves should be reduced to be able to fit the model.

#### ACKNOWLEDGEMENTS

The authors thank Dr S.E. Sallan and Dr L.B. Silverman for permission to use the acute lymphoblastic leukemia (ALL) data set. We would also like to thank V.K. Dalton and Dr M. Rue for helpful conversations regarding the scientific questions of interest in adolescent growth studies. The work of the first author was supported by DGES project BEC 2001-1270, the second author was supported by NIEHS grant ES07142, NIH grant GM29745, and Grant CA 68484 from the National Cancer Institute, and the research of the fourth author was supported by a grant from the National Cancer Institute (CA57030), and by the Texas A&M Center for Environmental and Rural Health via a grant from the National Institute of Environmental Health Sciences (P30-ES09106).

#### REFERENCES

1. N.M. Laird and J.H. Ware. Random-effects models for longitudinal data. *Biometrics*, 38:963–974, 1982.
2. S.L. Zeger and P. J. Diggle. Semiparametric models for longitudinal data with application to cd4 cell numebers in hiv seroconverters. *Biometrics*, 50:689–699, 1994.
3. A.P. Verbyla, B.R. Cullis, M.G. Kenward, and S.J. Welham. The analysis of designed experiments and longitudinal data using smoothing splines. *Applied Statistics*, 48:269–312, 1999.
4. D. Zhang, X. Lin, J. Raz, and M. Sowers. Semiparametric stochastic mixed models for longitudinal data. *Journal of the American Statistical Association*, 93:710–719, 1998.
5. B.A. Brumback and J.A. Rice. Smoothing spline models for the analysis of nested and crossed samples of curves. *Journal of the American Statistical Association*, 93:961–994, 1998.
6. J.A. Rice and C. O. Wu. Nonparametric mixed effects models for unequally sampled noisy curves. *Biometrics*, 57:253–259, 2001.

7. W. Guo. Functional mixed effects models. *Biometrics*, 58:121–128, 2002.
8. S.J. Koopman and J. Durbin. Fast filtering and smoothing for multivariate state space models. *Journal of Time Series Analysis*, 21:281–296, 2000.
9. P.H.C. Eilers and B.D. Marx. Flexible smoothing with  $B$ -splines and penalties. *Statistical Science*, 11:89–121, 1996.
10. B.A. Brumback, D. Ruppert, and M.P. Wand. Comment on “Variable selection and function estimation in additive nonparametric regression using a data-based prior”. *Journal of the American Statistical Association*, 94:794–797, 1999.
11. I. D. Currie and M. Durbán. Flexible smoothing with  $P$ -splines: A unified approach. *Statistical Modelling*, 2:333–349, 2002.
12. M.P. Wand. Smoothing and mixed models. *Computational Statistics*, 18:223–249, 2003.
13. V.K. Dalton, M. Rue, L.B. Silverman, R.D. Gelber, B.L. Asselin, R.D. Barr, L.A. Clavell, C. A. Hurwitz, A. Moghrabi, Samson Y., M. Schorin, N.J. Tarbell, S.E. Sallan, and L.E. Cohen. Height and weight in children treated for acute lymphoblastic leukemia: Relationship to cns treatment. *Journal of Clinical Oncology*, 21:2953–2960, 2003.
14. Y. Wang. Mixed effects smoothing spline analysis of variance. *Journal of the Royal Statistical Society B*, 60:159–174, 1998.
15. D. Ruppert, M.P. Wand, and R. J. Carroll. *Semiparametric Regression*. Cambridge University Press, 2003.
16. B.A. Coull, J. Schwartz, and M.P. Wand. Respiratory health and air pollution: Additive mixed model analyses. *Biostatistics*, 2:337–349, 2001.
17. H.D. Patterson and R. Thompson. Recovery of inter-block information when block sizes are unequal. *Biometrika*, 58:545–554, 1971.
18. S.G. Self and K. Liang. Asymptotic properties of maximum likelihood estimators and likelihood ratio tests under nonstandard conditions. *Journal of the American Statistical Association*, 82:605–610, 1987.
19. D.O. Stram and J.W. Lee. Variance components testing in the longitudinal mixed effects model. *Biometrics*, 50:1171–1177, 1994.
20. C.M. Crainiceanu, D. Ruppert, and T.J. Vogelsang. Probability that the mle of a variance component is zero with applications to likelihood ratio tests. *submitted*, 2002.
21. C.M. Crainiceanu and D. Ruppert. Likelihood ratio tests in linear mixed models with one variance component. *Journal of the Royal Statistical Society, B.*, 66:165–185, 2003.

22. C.M. Crainiceanu, D. Ruppert, G. Claeskens, and M.P. Wand. Exact likelihood ratio tests for penalized splines. *Biometrika*, to appear, 2004.
23. C.M. Crainiceanu and D. Ruppert. Restricted likelihood ratio tests for longitudinal models. *Statistica Sinica*, to appear, 2004.
24. A.C.S. Hokken-Koelega, J.W.D. van Doorn, K. Hahlen, T. Stijnen, S.M. de Muink Keizer-Schrama, and S.L. Drop. Long-term effects of treatment for acute lymphoblastic leukemia with and without cranial irradiation on growth and puberty. *Pediatric Research*, 33:577–582, 1993.
25. E. Didcock, D.A. Davies, M. Didi, A.L. Olgilvy Stuart, J.K. Wales, and S.M. Shalet. Pubertal growth in young adult survivors of childhood leukemia. *Journal of Clinical Oncology*, 1995.
26. D. Ruppert. Selecting the number of knots for penalized splines. *Journal of Computational and Graphical Statistics*, 11:735–757, 2002.
27. A.R. Gilmour, B.R. Cullis, S.J. Wellham, and R. Thompson. *ASREML Reference Manual*. NSW Agriculture, 2000.