

A bivariate semiparametric stochastic mixed model

Kexin Ji & Joel A. Dubin

August 22, 2017

1 Introduction

Longitudinal data analysis has wide applications in areas such as medicine, agriculture and so on. One distinctive feature of longitudinal data is that measurements of each subject are collected repeatedly over time; thus, each measurement for the same subject is no longer independent. This induces a correlation structure that needs to be modeled carefully. Many methods have been developed over the years to accommodate this additional structure.

Brumback and Rice [1] used natural cubic splines to model the mean structure of the linear mixed model. They extended the traditional LME model to generalized smoothing spline models for samples of curves stratified by nested and crossed factors and specified the design matrices associated with fixed effects and random effects by bases of functions, as opposed to the usual known covariates matrices. Verbyla et al [11] advocated a similar approach as Brumback & Rice [1], where data-based determination of the smoothing parameters was advocated in the paper, yet their model specification is slightly different; and the techniques were applied to the analysis of designed experiments.

In addition to modeling the mean structure of using a smoothing spline, some efforts were geared toward modeling complicated within-subject covariance. Taylor et al [10] used a particular stochastic process to model the data in addition to the usual random effect term which induces within-subject correlation. Zhang et al [16] combine both the smoothing spline to measure mean structure and various stationary and nonstationary stochastic processes to model serial correlation into one cohesive model. To further model cyclic responses, Zhang et al [15] extended their previous work and proposed a semiparametric stochastic mixed model for periodic longitudinal data. They used parametric functions to model the covariate effects and a periodic smooth nonparametric function to model the underlying complex periodic time course. The within-subject covariance is modeled using a random intercept and a stochastic process with periodic variance function. Instead of cubic smoothing spline, Welham et al [13] modelled cyclic longitudinal data using mixed model L-splines. Meyer et al [8] proposed a functional data analysis approach to model cyclic data. Last but not least, Wood [14] modified penalized cubic regression spline to model a cyclic smooth function.

All of the approaches mentioned above are to be applied on univariate longitudinal responses. A growing number of data require techniques to model bivariate, and in more generality, multivariate responses. Liu et al [7] extended the univariate state space model in time series analysis, and proposed a bivariate hierarchical state space model to bivariate longitudinal responses. Each response is modelled by a hierarchical state space model, with both population-average and subject-specific components. The bivariate model is constructed by linking the univariate models based on the hypothesized relationship. Sy, Taylor, and Cumberland [9] employed multivariate stochastic processes to jointly model bivariate longitudinal data.

We extended Zhang et al [16] [15] and propose a bivariate semiparametric stochastic mixed model for bivariate periodic repeated measures data. The bivariate model uses parametric fixed effects for modeling covariate effects and periodic smooth nonparametric functions for each of the two underlying time effects. In addition, the between-subject and within-subject correlations are modeled using separate but correlated random effects and a bivariate Gaussian random field, respectively. We derive maximum penalized likelihood estimators for both the fixed effects regression coefficients and the nonparametric time functions. The smoothing parameters and all variance components are estimated simultaneously using restricted maximum likelihood.

The paper is organized as followed. Section 2 specified the proposed model with assumptions. Section 3 provides estimation and inference procedures. Specifically, Sections 3.1 and 3.2 gives estimation procedures

for the model parameters, the nonparametric components, random effects and the Gaussian fields. Section 3.3 specifies the biases and covariances for all the estimators given in Section 3.2; and Section 3.4 concluded this section by providing the estimation procedures of the smoothing parameters and variance components. Section 4 investigate the proposed methodology through simulation. Section 5 illustrates the model by analyzing bivariate longitudinal female hormone data collected daily over a menstrual cycle. Section 6 discusses challenges and future work.

2 The Bivariate semiparametric stochastic mixed model

We propose a semiparametric stochastic bivariate mixed model, where the joint model assumes a semi-parametric mixed model for each outcome. The univariate models for each outcome are connected through the specification of the correlation structure for the random effects.

2.1 General bivariate model with joint distribution of random effects

Denote $\{Y_{1ij}, Y_{2ij}\}$ to be the bivariate response for the i th subject at time point j , $i = 1, \dots, m$ and $j = 1, \dots, n_i$. The bivariate model is written as

$$\begin{aligned} Y_{1ij} &= \mathbf{X}_{1ij}^T \boldsymbol{\beta}_1 + f_1(t_{ij}) + \mathbf{Z}_{1ij}^T \mathbf{b}_{1i} + U_{1i}(t_{ij}) + \epsilon_{1ij} \\ Y_{2ij} &= \mathbf{X}_{2ij}^T \boldsymbol{\beta}_2 + f_2(t_{ij}) + \mathbf{Z}_{2ij}^T \mathbf{b}_{2i} + U_{2i}(t_{ij}) + \epsilon_{2ij}, \end{aligned} \quad (1)$$

where $(\mathbf{X}_{1ij}, \mathbf{X}_{2ij})$ are known covariates associated with the fixed effects; $(\mathbf{Z}_{1ij}, \mathbf{Z}_{2ij})$ are known covariates associated with the random effects; $\boldsymbol{\beta}_1$ and $\boldsymbol{\beta}_2$ are $p_1 \times 1$ and $p_2 \times 1$ vectors of regression coefficients, containing the fixed effects, respectively; \mathbf{b}_{1i} and \mathbf{b}_{2i} are $q_1 \times 1$ and $q_2 \times 1$ vectors of random effects, respectively; $f_1(t)$ and $f_2(t)$ are twice-differentiable periodic smooth functions of time with periods T_1 and T_2 , respectively; $\{(U_{1i}(t_{ij}), U_{2i}(t_{ij})), t_{ij} \in \{t_{i1}, \dots, t_{in_i}\}, i = 1, \dots, m, j = 1, \dots, n_i\}$ are mean zero bivariate Gaussian field with covariance matrix

$$\begin{aligned} \mathbf{C}_i(s, t) &= \begin{pmatrix} E[U_{1i}(s)U_{1i}(t)] & E[U_{1i}(s)U_{2i}(t)] \\ E[U_{2i}(s)U_{1i}(t)] & E[U_{2i}(s)U_{2i}(t)] \end{pmatrix} \\ &= \begin{pmatrix} \sqrt{\xi_1(s)\xi_1(t)}\eta_1(\rho_1; s, t) & \sqrt{\xi_1(s)\xi_2(t)}\eta_3(\rho_3; s, t) \\ \sqrt{\xi_2(s)\xi_1(t)}\eta_3(\rho_3; t, s) & \sqrt{\xi_2(s)\xi_2(t)}\eta_2(\rho_2; s, t) \end{pmatrix} \end{aligned}$$

where $\xi_1(t)$ and $\xi_2(t)$ are periodic variance functions; $\text{corr}(U_{1i}(t), U_{1i}(s)) = \eta_1(\rho_1; s, t)$, $\text{corr}(U_{2i}(t), U_{2i}(s)) = \eta_2(\rho_2; s, t)$, and $\text{corr}(U_{1i}(t), U_{2i}(s)) = \eta_3(\rho_3; s, t)$ are correlation functions, where $\rho_1 \in [0, 1]$, $\rho_2 \in [0, 1]$ and $\rho_3 \in [0, 1]$ are correlation parameters; and the measurement errors $(\epsilon_{1ij}, \epsilon_{2ij})^T$ are bivariate normal

$$\begin{pmatrix} \epsilon_{1ij} \\ \epsilon_{2ij} \end{pmatrix} \sim N_2 \left(\mathbf{0}, \begin{pmatrix} \sigma_1^2 & \sigma_{12}^2 \\ \sigma_{12}^2 & \sigma_2^2 \end{pmatrix} \right).$$

We assume that \mathbf{b}_{ki} , $k = 1, 2$ to be $(q_1 + q_2)$ -dimensional normal with mean zero and covariance matrix $\mathbf{D}(\phi)$. These random effects, \mathbf{b}_{1i} and \mathbf{b}_{2i} , are assumed to be separate but correlated. Further, we assume that the random effects, the stochastic process and the measurement error to be mutually independent.

Denote

$$\mathbf{Y}_{ij} := \begin{pmatrix} Y_{1ij} \\ Y_{2ij} \end{pmatrix} \in \mathbb{R}^2$$

to be the response vector;

$$\mathbf{X}_{ij} := \begin{pmatrix} \mathbf{X}_{1ij}^T & \mathbf{0} \\ \mathbf{0} & \mathbf{X}_{2ij}^T \end{pmatrix} \in \mathbb{R}^{(p_1 + p_2) \times 2}, \quad \boldsymbol{\beta} := \begin{pmatrix} \boldsymbol{\beta}_1 \\ \boldsymbol{\beta}_2 \end{pmatrix} \in \mathbb{R}^{(p_1 + p_2)},$$

to be the matrix of known covariates and the vector of regression coefficients respectively;

$$\mathbf{Z}_{ij} := \begin{pmatrix} \mathbf{Z}_{1ij}^T & \mathbf{0} \\ \mathbf{0} & \mathbf{Z}_{2ij}^T \end{pmatrix} \in \mathbb{R}^{(q_1 + q_2) \times 2}, \quad \mathbf{b}_i := \begin{pmatrix} \mathbf{b}_{1i} \\ \mathbf{b}_{2i} \end{pmatrix} \in \mathbb{R}^{(q_1 + q_2)},$$

to be the matrix of known covariates associated with the random effects and the vector of random effects respectively; and finally

$$\mathbf{f}(t_{ij}) := \begin{pmatrix} f_1(t_{ij}) \\ f_2(t_{ij}) \end{pmatrix} \in \mathbb{R}^2, \quad \mathbf{U}_i(t_{ij}) := \begin{pmatrix} U_{1i}(t_{ij}) \\ U_{2i}(t_{ij}) \end{pmatrix} \in \mathbb{R}^2, \quad \boldsymbol{\epsilon}_{ij} := \begin{pmatrix} \epsilon_{1ij} \\ \epsilon_{2ij} \end{pmatrix} \in \mathbb{R}^2,$$

to be the vectors of the smooth function, the stochastic process, and the measurement error of. Then model (1) can also be rewritten as

$$\mathbf{Y}_{ij} = \mathbf{X}_{ij}^T \boldsymbol{\beta} + \mathbf{f}(t_{ij}) + \mathbf{Z}_{ij}^T \mathbf{b}_i + \mathbf{U}_i(t_{ij}) + \boldsymbol{\epsilon}_{ij}, \quad (2)$$

with the same model assumptions.

This model (1) is an extension to the model proposed in Zhang et al [16] [15], where a univariate semiparametric stochastic mixed model for (periodic) longitudinal data was proposed. The challenge here is that we are modeling a bivariate longitudinal response model, which is achieved by modeling a joint distribution of the random effects. The distributions of the two random effects can be potentially distinct, with different distributions or the same distribution with different parameters; but, as mentioned above, the two random effects are assumed separate but correlated.

2.2 The Gaussian Field Specification

To accommodate for more complicated within-subject correlation, we propose to include various stationary and nonstationary Gaussian field to model serial correlation. This allows for variance to be varied over time.

There are potentially many choices available: Wiener process or Brownian motion (Taylor et al [10]); an integrated Wiener process and so on. One particular Gaussian process/field worthy of mentioning is Ornstein-Uhlenbeck (OU) process [6] which has a correlation function that decays exponentially over time $\text{corr}(U_i(t), U_i(s)) = \exp\{-\alpha|s - t|\}$. The variance function for OU process $\xi(t) = \sigma^2/2a$ is a constant, thus the process is strictly stationary. When $\xi(t)$ varies over time, then the process become nonhomogeneous (NOU) and for example we can assume $\xi(t) = \exp(a_0 + a_1 t)$.

3 Estimation and Inference

3.1 Matrix notation

To make inference from the model (2), we will write the model in matrix form - first, in subject level; then, over all subjects. Denote

$$\mathbf{Y}_i := \begin{pmatrix} \mathbf{Y}_{i1} \\ \vdots \\ \mathbf{Y}_{in_i} \end{pmatrix} \in \mathbb{R}^{2n_i},$$

to be the response vector;

$$\mathbf{X}_i := \begin{pmatrix} \mathbf{X}_{i1}^T \\ \vdots \\ \mathbf{X}_{in_i}^T \end{pmatrix} \in \mathbb{R}^{2n_i \times (p_1 + p_2)}, \quad \mathbf{Z}_i := \begin{pmatrix} \mathbf{Z}_{i1}^T \\ \vdots \\ \mathbf{Z}_{in_i}^T \end{pmatrix} \in \mathbb{R}^{2n_i \times (q_1 + q_2)},$$

to be the corresponding covariate matrix associate with the fixed effects and the random effects respectively; and

$$\mathbf{U}_i := \begin{pmatrix} \mathbf{U}_i(t_{i1}) \\ \vdots \\ \mathbf{U}_i(t_{in_i}) \end{pmatrix} \in \mathbb{R}^{2n_i}, \quad \boldsymbol{\epsilon}_i := \begin{pmatrix} \boldsymbol{\epsilon}_{i1} \\ \vdots \\ \boldsymbol{\epsilon}_{in_i} \end{pmatrix} \in \mathbb{R}^{2n_i},$$

to be the vectors of stochastic process and measurement errors. Assume $t_{ij} > 0$ and $\min\{t_{ij}\} = 0$. Since $f_1(t)$ and $f_2(t)$ are periodic functions with periods T_1 and T_2 , we only need to estimate $f_1(t)$ for $t \in [0, T_1)$ and

$f_2(t)$ for $t \in [0, T_2)$. Let $\mathbf{t}'_1 = (t'_{11}, \dots, t'_{1r_1})$ to be a vector of ordered distinct values of $t'_{1ij} = \text{mod}(t_{ij}, T_1)$ for $i = 1, \dots, m$ and $j = 1 \dots n_i$, and let $\mathbf{t}'_2 = (t'_{21}, \dots, t'_{2r_2})$ to be a vector of ordered distinct values of $t'_{2ij} = \text{mod}(t_{ij}, T_2)$ for $i = 1, \dots, m$ and $j = 1 \dots n_i$, thus $t'_{1k} \in [0, T_1)$ for $k = 1, \dots, r_1$ and $t'_{2k} \in [0, T_2)$ for $k = 1, \dots, r_2$. Then, let $\tilde{\mathbf{N}}_{1i}$ be the $n_i \times r_1$ incidence matrix for the i^{th} subject for the first response connecting $\mathbf{t}_i = (t_{i1}, \dots, t_{in_i})^T$ and \mathbf{t}'_1 such that

$$\tilde{\mathbf{N}}_{1i}[j, \ell] = \begin{cases} 1 & \text{if } t'_{1ij} = t'_{1\ell} \\ 0 & \text{otherwise,} \end{cases}$$

where $\tilde{\mathbf{N}}_{1i}[j, \ell]$ denote the $(j, \ell)^{\text{th}}$ entry of matrix $\tilde{\mathbf{N}}_{1i}$ for $j = 1, \dots, n_i$ and $\ell = 1, \dots, r_1$. Similarly, let $\tilde{\mathbf{N}}_{2i}$ be the $n_i \times r_2$ incidence matrix for the i^{th} subject for the second response connecting $\mathbf{t}_i = (t_{i1}, \dots, t_{in_i})^T$ and \mathbf{t}'_2 such that

$$\tilde{\mathbf{N}}_{2i}[j, \ell] = \begin{cases} 1 & \text{if } t'_{2ij} = t'_{2\ell} \\ 0 & \text{otherwise,} \end{cases}$$

where $\tilde{\mathbf{N}}_{2i}[j, \ell]$ denote the $(j, \ell)^{\text{th}}$ entry of matrix $\tilde{\mathbf{N}}_{2i}$ for $j = 1, \dots, n_i$ and $\ell = 1, \dots, r_2$. Further, we need to refine the incidence matrix $\tilde{\mathbf{N}}_{1i}$ to make it correspond to the first response such that

$$\mathbf{N}_{1i} = \mathbf{A}_{1i} \tilde{\mathbf{N}}_{1i}$$

where

$$\mathbf{A}_{1i} = \begin{pmatrix} 1 & 0 & 0 & \dots & 0 \\ 0 & 0 & 0 & \dots & 0 \\ 0 & 1 & 0 & \dots & 0 \\ 0 & 0 & 0 & \dots & 0 \\ \vdots & & & \ddots & \\ 0 & \dots & 0 & 0 & 1 \\ 0 & 0 & 0 & \dots & 0 \end{pmatrix} \in \mathbb{R}^{2n_i \times n_i},$$

thus the refined incidence matrix \mathbf{N}_{1i} is of dimension $2n_i \times r_1$. Similarly, the refined incidence matrix \mathbf{N}_{2i} of dimension $2n_i \times r_2$ for the second response is

$$\mathbf{N}_{2i} = \mathbf{A}_{2i} \tilde{\mathbf{N}}_{2i}$$

where

$$\mathbf{A}_{2i} = \begin{pmatrix} 0 & 0 & 0 & \dots & 0 \\ 1 & 0 & 0 & \dots & 0 \\ 0 & 0 & 0 & \dots & 0 \\ 0 & 1 & 0 & \dots & 0 \\ \vdots & & & \ddots & \\ 0 & 0 & 0 & \dots & 0 \\ 0 & \dots & 0 & 0 & 1 \end{pmatrix} \in \mathbb{R}^{2n_i \times n_i}.$$

Then the proposed bivariate semiparametric stochastic mixed model (1) can be written as

$$\mathbf{Y}_i = \mathbf{X}_i \boldsymbol{\beta} + \mathbf{N}_{1i} \mathbf{f}_1 + \mathbf{N}_{2i} \mathbf{f}_2 + \mathbf{Z}_i \mathbf{b}_i + \mathbf{U}_i + \boldsymbol{\epsilon}_i$$

for subject i , where

$$\mathbf{f}_1 := \begin{pmatrix} f_1(t'_{11}) \\ \vdots \\ f_1(t'_{1r_1}) \end{pmatrix} \in \mathbb{R}^{r_1}, \quad \mathbf{f}_2 := \begin{pmatrix} f_2(t'_{21}) \\ \vdots \\ f_2(t'_{2r_2}) \end{pmatrix} \in \mathbb{R}^{r_2}.$$

Further denoting $\mathbf{Y} = (\mathbf{Y}_1^T, \dots, \mathbf{Y}_m^T)^T$ and $\mathbf{X}, \mathbf{N}_1, \mathbf{N}_2, \mathbf{b}, \mathbf{U}, \boldsymbol{\epsilon}$ similarly and let $n = \sum_{i=1}^m n_i$, then the bivariate semiparametric stochastic mixed effects model over all subjects is written as

$$\mathbf{Y} = \mathbf{X} \boldsymbol{\beta} + \mathbf{N}_1 \mathbf{f}_1 + \mathbf{N}_2 \mathbf{f}_2 + \mathbf{Z} \mathbf{b} + \mathbf{U} + \boldsymbol{\epsilon} \quad (3)$$

where

$$\begin{pmatrix} b \\ U \\ \epsilon \end{pmatrix} \sim N \left(\begin{pmatrix} 0 \\ 0 \\ 0 \end{pmatrix}, \begin{pmatrix} D(\phi) & 0 & 0 \\ 0 & \Gamma(\xi, \rho) & 0 \\ 0 & 0 & \Sigma(\sigma^2) \end{pmatrix} \right)$$

with $D(\phi) = \text{diag}(D, \dots, D)$; $\Gamma(\xi, \rho) = \text{diag}(\Gamma_1(t_1, t_1), \dots, \Gamma_m(t_m, t_m))$ and the $(k, k')^{\text{th}}$ entry of $\Gamma_i(t_i, t_i)$ is $C_i(k, k')$; and $\Sigma(\sigma^2) = \begin{pmatrix} \sigma_1^2 & 0 \\ 0 & \sigma_2^2 \end{pmatrix} I_n$.

3.2 Estimation of Model Coefficients, Nonparametric Function, Random Effects and Gaussian Fields

Let $Y = X\beta + N_1 f_1 + N_2 f_2 + \epsilon^*$, where

$$\epsilon^* = Zb + U + \epsilon = \begin{pmatrix} Z & I_{2n \times 2n} & I_{2n \times 2n} \end{pmatrix} \begin{pmatrix} b \\ U \\ \epsilon \end{pmatrix}$$

Then $\epsilon^* \sim N_{2n}(\mathbf{0}, V)$ where

$$\begin{aligned} V &= \text{Acov} \begin{pmatrix} b \\ U \\ \epsilon \end{pmatrix} A^T = \begin{pmatrix} Z & I_{2n \times 2n} & I_{2n \times 2n} \end{pmatrix} \begin{pmatrix} D(\phi) & 0 & 0 \\ 0 & \Gamma(\xi, \rho) & 0 \\ 0 & 0 & \Sigma(\sigma) \end{pmatrix} \begin{pmatrix} Z^T \\ I_{2n \times 2n} \\ I_{2n \times 2n} \end{pmatrix} \\ &= \begin{pmatrix} ZD(\phi) & \Gamma(\xi, \rho) & \Sigma(\sigma) \end{pmatrix} \begin{pmatrix} Z^T \\ I_{2n \times 2n} \\ I_{2n \times 2n} \end{pmatrix} = ZDZ^T + \Gamma + \Sigma \end{aligned}$$

Therefore the proposed model (3) also implies the *marginal model*

$$Y = X\beta + N_1 f_1 + N_2 f_2 + \epsilon^*, \quad \epsilon^* \sim N_{2n}(\mathbf{0}, V)$$

where $V = ZDZ^T + \Gamma + \Sigma$. By the marginal model (4), the *log-likelihood* function for (β, f_1, f_2) :

$$\ell(\beta, f_1, f_2; Y) = -\frac{1}{2} \log |V| - \frac{1}{2} (Y - \beta - N_1 f_1 - N_2 f_2)^T V^{-1} (Y - \beta - N_1 f_1 - N_2 f_2)$$

We estimate the parameters β , f_1 and f_2 by maximizing the penalized likelihood [12]:

$$\ell(\beta, f_1, f_2; Y) - \lambda_1 \int_a^b [f_1''(t)]^2 dt - \lambda_2 \int_a^b [f_2''(t)]^2 dt = \ell(\beta, f_1, f_2; Y) - \lambda_1 f_1^T K f_1 - \lambda_2 f_2^T K f_2 \quad (4)$$

where K is the nonnegative definite smoothing matrix, defined in Equation (2.3) in Green and Silverman [5]. And the resulting estimators for the nonparametric functions are the natural cubic spline estimators of f_1 and f_2 .

Given fixed smoothing parameters and variance parameters, differentiation of (4) with respect to β , f_1 , f_2 gives the estimators $(\hat{\beta}, \hat{f}_1, \hat{f}_2)$ that solves

$$\begin{pmatrix} X^T W X & X^T W N_1 & X^T W N_2 \\ N_1^T W X & N_1^T W N_1 + \lambda_1 K & N_1^T W N_2 \\ N_2^T W X & N_2^T W N_1 & N_2^T W N_2 + \lambda_2 K \end{pmatrix} \begin{pmatrix} \beta \\ f_1 \\ f_2 \end{pmatrix} = \begin{pmatrix} X^T W Y \\ N_1^T W Y \\ N_2^T W Y \end{pmatrix},$$

where $W = V^{-1}$. To study the theoretical properties of the estimates, such as bias and covariance, we derive the closed-form solutions for $\hat{\beta}$, \hat{f}_1 and \hat{f}_2

$$\hat{\beta} = (X^T W_x X)^{-1} X^T W_x Y \quad (5)$$

$$\hat{f}_1 = (N_1^T W_{f_1} N_1 + \lambda_1 K)^{-1} N_1^T W_{f_1} Y \quad (6)$$

$$\hat{f}_2 = (N_2^T W_{f_2} N_2 + \lambda_2 K)^{-1} N_2^T W_{f_2} Y, \quad (7)$$

where $\mathbf{W}_x = \mathbf{W}_1 - \mathbf{W}_1 \mathbf{N}_2 (\mathbf{N}_2^T \mathbf{W}_1 \mathbf{N}_2 + \lambda_2 \mathbf{K})^{-1} \mathbf{N}_2^T \mathbf{W}_1$, $\mathbf{W}_{f_1} = \mathbf{W}_2 - \mathbf{W}_2 \mathbf{X} (\mathbf{X}^T \mathbf{W}_2 \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W}_2$, and $\mathbf{W}_{f_2} = \mathbf{W}_1 - \mathbf{W}_1 \mathbf{X} (\mathbf{X}^T \mathbf{W}_1 \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W}_1$ are weight matrices with $\mathbf{W}_1 = \mathbf{W} - \mathbf{W} \mathbf{N}_1 (\mathbf{N}_1^T \mathbf{W} \mathbf{N}_1 + \lambda_1 \mathbf{K})^{-1} \mathbf{N}_1^T \mathbf{W}$ and $\mathbf{W}_2 = \mathbf{W} - \mathbf{W} \mathbf{N}_2 (\mathbf{N}_2^T \mathbf{W} \mathbf{N}_2 + \lambda_2 \mathbf{K})^{-1} \mathbf{N}_2^T \mathbf{W}$.

Estimation of the subject-specific random effects \mathbf{b}_i and the subject-specific Gaussian field $\mathbf{U}_i(\mathbf{s}_i)$ is obtained by calculating their conditional expectations given the data \mathbf{Y} . Note that the proposed model (3) can also be rewritten as *two-level hierarchical model*

$$\mathbf{Y} | \mathbf{b}, \mathbf{U} \sim N_{2n}(\mathbf{X}\boldsymbol{\beta} + \mathbf{N}_1 \mathbf{f}_1 + \mathbf{N}_2 \mathbf{f}_2 + \mathbf{Z}\mathbf{b} + \mathbf{U}, \boldsymbol{\Sigma}) \quad (8)$$

$$\mathbf{b} \sim N_{2m}(\mathbf{0}, \mathbf{D}) \quad (9)$$

$$\mathbf{U} \sim N_{2n}(\mathbf{0}, \boldsymbol{\Gamma}),$$

then by the property of Normality, we have

$$\begin{pmatrix} \mathbf{Y} \\ \mathbf{b} \end{pmatrix} \sim N_{2n+2m} \left(\begin{pmatrix} \mathbf{X}\boldsymbol{\beta} + \mathbf{N}_1 \mathbf{f}_1 + \mathbf{N}_2 \mathbf{f}_2 \\ \mathbf{0} \end{pmatrix}, \begin{pmatrix} \mathbf{V} & \mathbf{Z}\mathbf{D} \\ \mathbf{D}\mathbf{Z}^T & \mathbf{D} \end{pmatrix} \right).$$

since

$$\begin{aligned} \text{Cov}(\mathbf{Y}, \mathbf{b}) &= \text{Cov}(\mathbf{X}\boldsymbol{\beta} + \mathbf{N}_1 \mathbf{f}_1 + \mathbf{N}_2 \mathbf{f}_2 + \mathbf{Z}\mathbf{b} + \mathbf{U} + \boldsymbol{\epsilon}, \mathbf{b}) \\ &= \text{Cov}(\mathbf{X}\boldsymbol{\beta}, \mathbf{b}) + \text{Cov}(\mathbf{N}_1 \mathbf{f}_1, \mathbf{b}) + \text{Cov}(\mathbf{N}_2 \mathbf{f}_2, \mathbf{b}) + \mathbf{Z}\text{Cov}(\mathbf{b}, \mathbf{b}) + \text{Cov}(\mathbf{U}, \mathbf{b}) + \text{Cov}(\boldsymbol{\epsilon}, \mathbf{b}) \\ &= \mathbf{Z}\mathbf{D} \end{aligned}$$

Therefore,

$$E(\mathbf{b} | \mathbf{Y}) = \mathbf{0} + \mathbf{D}\mathbf{Z}^T \mathbf{V}^{-1}(\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}} - \mathbf{N}_1 \hat{\mathbf{f}}_1 - \mathbf{N}_2 \hat{\mathbf{f}}_2) = \mathbf{D}\mathbf{Z}^T \mathbf{V}^{-1}(\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}} - \mathbf{N}_1 \hat{\mathbf{f}}_1 - \mathbf{N}_2 \hat{\mathbf{f}}_2)$$

and the estimator or predictor for subject-specific random effects \mathbf{b}_i is

$$\hat{\mathbf{b}}_i = \mathbf{D}\mathbf{Z}_i^T \mathbf{V}_i^{-1}(\mathbf{Y}_i - \mathbf{X}_i \hat{\boldsymbol{\beta}} - \hat{\mathbf{f}}_{1i} - \hat{\mathbf{f}}_{2i}) \quad (10)$$

Similarly, the estimator or predictor for the subject-specific Gaussian field $\mathbf{U}_i(\mathbf{s}_i)$ is

$$\hat{\mathbf{U}}_i(\mathbf{s}_i) = \boldsymbol{\Gamma}(\mathbf{s}_i, \mathbf{t}_i) \mathbf{V}_i^{-1}(\mathbf{Y}_i - \mathbf{X}_i \hat{\boldsymbol{\beta}} - \hat{\mathbf{f}}_{1i} - \hat{\mathbf{f}}_{2i}) \quad (11)$$

where $\hat{\mathbf{f}}_{1i} = \mathbf{N}_{1i} \hat{\mathbf{f}}_1$ and $\hat{\mathbf{f}}_{2i} = \mathbf{N}_{2i} \hat{\mathbf{f}}_2$.

3.3 Biases and Covariances of Model Coefficients, Nonparametric Function, Random Effects and Gaussian Fields

From closed-form solutions of estimators from equation (5) (6) and (7) in Section 3.2, the biases of the estimators $\hat{\boldsymbol{\beta}}$, $\hat{\mathbf{f}}_1$ and $\hat{\mathbf{f}}_2$ can be easily calculated and we have

$$E(\hat{\boldsymbol{\beta}}) - \boldsymbol{\beta} = (\mathbf{X}^T \mathbf{W}_x \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W}_x (\mathbf{N}_1 \mathbf{f}_1 + \mathbf{N}_2 \mathbf{f}_2)$$

$$E(\hat{\mathbf{f}}_1) - \mathbf{f}_1 = (\mathbf{N}_1^T \mathbf{W}_{f_1} \mathbf{N}_1 + \lambda_1 \mathbf{K})^{-1} (\mathbf{N}_1^T \mathbf{W}_{f_1} \mathbf{N}_2 \mathbf{f}_2 - \lambda_1 \mathbf{K} \mathbf{f}_1),$$

and

$$E(\hat{\mathbf{f}}_2) - \mathbf{f}_2 = (\mathbf{N}_2^T \mathbf{W}_{f_2} \mathbf{N}_2 + \lambda_2 \mathbf{K})^{-1} (\mathbf{N}_2^T \mathbf{W}_{f_2} \mathbf{N}_1 \mathbf{f}_1 - \lambda_2 \mathbf{K} \mathbf{f}_2).$$

Similarly, the expected values of the estimators in (10) and (11) for the subject-specific random effects \mathbf{b}_i and for the subject-specific Gaussian field $\mathbf{U}_i(\mathbf{s}_i)$ are

$$\begin{aligned} E(\hat{\mathbf{b}}_i) &= \mathbf{D}\mathbf{Z}_i^T \mathbf{W}_i [\lambda_1 \mathbf{N}_{1i} (\mathbf{N}_1^T \mathbf{W}_{f_1} \mathbf{N}_1 + \lambda_1 \mathbf{K})^{-1} \mathbf{K} - \mathbf{X}_i (\mathbf{X}^T \mathbf{W}_x \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W}_x \mathbf{N}_1 \\ &\quad - \mathbf{N}_{2i} (\mathbf{N}_2^T \mathbf{W}_{f_2} \mathbf{N}_2 + \lambda_2 \mathbf{K})^{-1} \mathbf{N}_2^T \mathbf{W}_{f_2} \mathbf{N}_1] \mathbf{f}_1 \\ &\quad + \mathbf{D}\mathbf{Z}_i^T \mathbf{W}_i [\lambda_2 \mathbf{N}_{2i} (\mathbf{N}_2^T \mathbf{W}_{f_2} \mathbf{N}_2 + \lambda_2 \mathbf{K})^{-1} \mathbf{K} - \mathbf{X}_i (\mathbf{X}^T \mathbf{W}_x \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W}_x \mathbf{N}_2 \\ &\quad - \mathbf{N}_{1i} (\mathbf{N}_1^T \mathbf{W}_{f_1} \mathbf{N}_1 + \lambda_1 \mathbf{K})^{-1} \mathbf{N}_1^T \mathbf{W}_{f_1} \mathbf{N}_2] \mathbf{f}_2 \end{aligned}$$

and

$$\begin{aligned}
E \left[\hat{U}_i(\mathbf{s}_i) \right] &= \mathbf{\Gamma}_i(\mathbf{s}_i, \mathbf{t}_i) [\lambda_1 \mathbf{N}_{1i} (\mathbf{N}_1^T \mathbf{W}_{f_1} \mathbf{N}_1 + \lambda_1 \mathbf{K})^{-1} \mathbf{K} - \mathbf{X}_i (\mathbf{X}^T \mathbf{W}_x \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W}_x \mathbf{N}_1 \\
&\quad - \mathbf{N}_{2i} (\mathbf{N}_2^T \mathbf{W}_{f_2} \mathbf{N}_2 + \lambda_2 \mathbf{K})^{-1} \mathbf{N}_2^T \mathbf{W}_{f_2} \mathbf{N}_1] \mathbf{f}_1 \\
&\quad + \mathbf{\Gamma}_i(\mathbf{s}_i, \mathbf{t}_i) [\lambda_2 \mathbf{N}_{2i} (\mathbf{N}_2^T \mathbf{W}_{f_2} \mathbf{N}_2 + \lambda_2 \mathbf{K})^{-1} \mathbf{K} - \mathbf{X}_i (\mathbf{X}^T \mathbf{W}_x \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W}_x \mathbf{N}_2 \\
&\quad - \mathbf{N}_{1i} (\mathbf{N}_1^T \mathbf{W}_{f_1} \mathbf{N}_1 + \lambda_1 \mathbf{K})^{-1} \mathbf{N}_1^T \mathbf{W}_{f_1} \mathbf{N}_2] \mathbf{f}_2.
\end{aligned}$$

It can be shown that the biases of $\hat{\beta}$, $\hat{\mathbf{f}}_1$, $\hat{\mathbf{f}}_2$, $\hat{\mathbf{b}}_i$ and \hat{U}_i all go to $\mathbf{0}$ as $\lambda_1 \rightarrow 0$ and $\lambda_2 \rightarrow 0$.

For covariances, simple calculation using (5) (6) and (7) give the covariance of $\hat{\beta}$

$$\text{Cov}(\hat{\beta}) = (\mathbf{X}^T \mathbf{W}_x \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W}_x \mathbf{V} \mathbf{W}_x \mathbf{X} (\mathbf{X}^T \mathbf{W}_x \mathbf{X})^{-1}$$

and the covariance of $\hat{\mathbf{f}}_1$ and $\hat{\mathbf{f}}_2$

$$\text{Cov}(\hat{\mathbf{f}}_1) = (\mathbf{N}_1^T \mathbf{W}_{f_1} \mathbf{N}_1 + \lambda_1 \mathbf{K})^{-1} \mathbf{N}_1^T \mathbf{W}_{f_1} \mathbf{V} \mathbf{W}_{f_1} \mathbf{N}_1 (\mathbf{N}_1^T \mathbf{W}_{f_1} \mathbf{N}_1 + \lambda_1 \mathbf{K})^{-1}$$

$$\text{Cov}(\hat{\mathbf{f}}_2) = (\mathbf{N}_2^T \mathbf{W}_{f_2} \mathbf{N}_2 + \lambda_2 \mathbf{K})^{-1} \mathbf{N}_2^T \mathbf{W}_{f_2} \mathbf{V} \mathbf{W}_{f_2} \mathbf{N}_2 (\mathbf{N}_2^T \mathbf{W}_{f_2} \mathbf{N}_2 + \lambda_2 \mathbf{K})^{-1}.$$

The covariances of the estimators in (10) and (11) for the subject-specific random effects \mathbf{b}_i and for the subject-specific Gaussian field $U_i(\mathbf{s}_i)$ are

$$\text{Cov}(\hat{\mathbf{b}}_i - \mathbf{b}_i) = \mathbf{D} - \mathbf{D} \mathbf{Z}_i^T \mathbf{W}_i \mathbf{Z}_i \mathbf{D} + \mathbf{D} \mathbf{Z}_i^T \mathbf{W}_i \chi_i \mathbf{C}^{-1} \chi_i^T \mathbf{W}_i \mathbf{C}^{-1} \chi_i^T \mathbf{W}_i \mathbf{Z}_i \mathbf{D}$$

and

$$\text{Cov}(\hat{U}_i(\mathbf{s}_i) - U_i(\mathbf{s}_i)) = \mathbf{\Gamma}(\mathbf{s}_i, \mathbf{s}_i) - \mathbf{\Gamma}(\mathbf{s}_i, \mathbf{t}_i) \mathbf{W}_i \mathbf{\Gamma}(\mathbf{s}_i, \mathbf{t}_i)^T + \mathbf{\Gamma}(\mathbf{s}_i, \mathbf{t}_i) \mathbf{W}_i \chi_i \mathbf{C}^{-1} \chi_i^T \mathbf{W}_i \mathbf{C}^{-1} \chi_i^T \mathbf{W}_i \mathbf{\Gamma}(\mathbf{s}_i, \mathbf{t}_i)^T,$$

where $\chi_i = (\mathbf{X}_i \quad \mathbf{N}_{1i} \quad \mathbf{N}_{2i})$ and $\chi = (\mathbf{X} \quad \mathbf{N}_1 \quad \mathbf{N}_2)$.

3.4 Estimation of the Smoothing Parameters and Variance Parameters

3.4.1 The Linear Mixed Model Representation

By Green (1997) [4], we can write \mathbf{f}_1 and \mathbf{f}_2 by a one-to-one linear transformation as

$$\begin{aligned}
\mathbf{f}_1 &= \mathbf{T}_1 \delta_1 + \mathbf{B}_1 \mathbf{a}_1 \\
\mathbf{f}_2 &= \mathbf{T}_2 \delta_2 + \mathbf{B}_2 \mathbf{a}_2
\end{aligned}$$

where δ_1 and \mathbf{a}_1 are of dimensions 2 and $r_1 - 2$ and δ_2 and \mathbf{a}_2 are of dimensions 2 and $r_2 - 2$. $\mathbf{B}_1 = \mathbf{L}_1 (\mathbf{L}_1^T \mathbf{L}_1)^{-1}$ and \mathbf{L}_1 is $r_1 \times (r_1 - 2)$ full-rank matrix satisfying $\mathbf{K}_1 = \mathbf{L}_1 \mathbf{L}_1^T$ and $\mathbf{L}_1^T \mathbf{T}_1 = \mathbf{0}$. $\mathbf{B}_2 = \mathbf{L}_2 (\mathbf{L}_2^T \mathbf{L}_2)^{-1}$ and \mathbf{L}_2 is $r_2 \times (r_2 - 2)$ full-rank matrix satisfying $\mathbf{K}_2 = \mathbf{L}_2 \mathbf{L}_2^T$ and $\mathbf{L}_2^T \mathbf{T}_2 = \mathbf{0}$.

Thus the proposed semiparametric mixed model (3) can be rewritten as a modified linear mixed model,

$$\mathbf{Y} = \mathbf{X} \beta + \mathbf{N}_1 \mathbf{T}_1 \delta_1 + \mathbf{N}_1 \mathbf{B}_1 \mathbf{a}_1 + \mathbf{N}_2 \mathbf{T}_2 \delta_2 + \mathbf{N}_2 \mathbf{B}_2 \mathbf{a}_2 + \mathbf{Z} \mathbf{b} + \mathbf{U} + \epsilon, \quad (12)$$

where $\beta_* = (\beta^T, \delta_1^T, \delta_2^T)^T$ are the regression coefficients and $\mathbf{b}_* = (\mathbf{a}_1^T, \mathbf{a}_2^T, \mathbf{b}^T, \mathbf{U}^T)^T$ are mutually independent random effects with \mathbf{a}_1 distributed as normal $(0, \tau_1 \mathbf{I})$, \mathbf{a}_2 distributed as normal $(0, \tau_2 \mathbf{I})$, and (\mathbf{b}, \mathbf{U}) having the same distribution as specified before. The marginal variance of \mathbf{Y} under the modified mixed model representation becomes $\mathbf{V}_* = \tau_1 \mathbf{B}_{1*} \mathbf{B}_{1*}^T + \tau_2 \mathbf{B}_{2*} \mathbf{B}_{2*}^T + \mathbf{V}$, where $\mathbf{B}_{1*} = \mathbf{N}_1 \mathbf{B}_1$ and $\mathbf{B}_{2*} = \mathbf{N}_2 \mathbf{B}_2$.

3.4.2 The Linear Mixed Model Representation

Under the above modified linear mixed model (12), the REML log-likelihood of (τ_1, τ_2, θ) is

$$\begin{aligned}
\ell_R(\tau_1, \tau_2, \theta; \mathbf{Y}) &= -\frac{1}{2} \log |\mathbf{V}_*| - \frac{1}{2} \log |\mathbf{X}_*^T \mathbf{V}_*^{-1} \mathbf{X}_*| - \frac{1}{2} (\mathbf{Y} - \mathbf{X}_* \hat{\beta}_*)^T \mathbf{V}_*^{-1} (\mathbf{Y} - \mathbf{X}_* \hat{\beta}_*) \\
&= -\frac{1}{2} \left[\log |\mathbf{V}_*| + \log |\mathbf{X}_*^T \mathbf{V}_*^{-1} \mathbf{X}_*| + (\mathbf{Y} - \mathbf{X}_* \hat{\beta}_*)^T \mathbf{V}_*^{-1} (\mathbf{Y} - \mathbf{X}_* \hat{\beta}_*) \right]
\end{aligned}$$

where $\mathbf{X}_* = [\mathbf{X}, \mathbf{N}_1 \mathbf{T}_1, \mathbf{N}_2 \mathbf{T}_2]$. Taking derivative with respect to τ_1 , τ_2 , and $\boldsymbol{\theta}$ and using the identity

$$\mathbf{V}_*^{-1}(\mathbf{Y} - \mathbf{X}_* \hat{\boldsymbol{\beta}}_*) = \mathbf{V}^{-1}(\mathbf{Y} - \mathbf{X} \hat{\boldsymbol{\beta}} - \mathbf{N}_1 \hat{\mathbf{f}}_1 - \mathbf{N}_2 \hat{\mathbf{f}}_2),$$

the estimating equation for the smoothing parameters τ_1 τ_2 and variance components $\boldsymbol{\theta}$ can be obtained

$$\frac{\partial \ell_R}{\partial \tau_1} = -\frac{1}{2} \text{Tr}(\mathbf{P}_* \mathbf{B}_{1*} \mathbf{B}_{1*}^T) + \frac{1}{2} (\mathbf{Y} - \mathbf{X} \hat{\boldsymbol{\beta}} - \mathbf{N}_1 \hat{\mathbf{f}}_1 - \mathbf{N}_2 \hat{\mathbf{f}}_2)^T \mathbf{V}^{-1} \mathbf{B}_{1*} \mathbf{B}_{1*}^T \mathbf{V}^{-1} (\mathbf{Y} - \mathbf{X} \hat{\boldsymbol{\beta}} - \mathbf{N}_1 \hat{\mathbf{f}}_1 - \mathbf{N}_2 \hat{\mathbf{f}}_2), \quad (13)$$

$$\frac{\partial \ell_R}{\partial \tau_2} = -\frac{1}{2} \text{Tr}(\mathbf{P}_* \mathbf{B}_{2*} \mathbf{B}_{2*}^T) + \frac{1}{2} (\mathbf{Y} - \mathbf{X} \hat{\boldsymbol{\beta}} - \mathbf{N}_1 \hat{\mathbf{f}}_1 - \mathbf{N}_2 \hat{\mathbf{f}}_2)^T \mathbf{V}^{-1} \mathbf{B}_{2*} \mathbf{B}_{2*}^T \mathbf{V}^{-1} (\mathbf{Y} - \mathbf{X} \hat{\boldsymbol{\beta}} - \mathbf{N}_1 \hat{\mathbf{f}}_1 - \mathbf{N}_2 \hat{\mathbf{f}}_2), \quad (14)$$

and

$$\frac{\partial \ell_R}{\partial \theta_j} = -\frac{1}{2} \text{Tr}(\mathbf{P}_* \frac{\partial \mathbf{V}}{\partial \theta_j}) + \frac{1}{2} (\mathbf{Y} - \mathbf{X} \hat{\boldsymbol{\beta}} - \mathbf{N}_1 \hat{\mathbf{f}}_1 - \mathbf{N}_2 \hat{\mathbf{f}}_2)^T \mathbf{V}^{-1} \frac{\partial \mathbf{V}}{\partial \theta_j} \mathbf{V}^{-1} (\mathbf{Y} - \mathbf{X} \hat{\boldsymbol{\beta}} - \mathbf{N}_1 \hat{\mathbf{f}}_1 - \mathbf{N}_2 \hat{\mathbf{f}}_2), \quad (15)$$

where $\mathbf{P}_* = \mathbf{V}_*^{-1} - \mathbf{V}_*^{-1} \mathbf{X}_* (\mathbf{X}_*^T \mathbf{V}_*^{-1} \mathbf{X}_*)^{-1} \mathbf{X}_*^T \mathbf{V}_*^{-1}$ is the projection matrix.

The covariance of the the smoothing parameters τ_1 τ_2 and variance components $\boldsymbol{\theta}$ can be estimated using Fisher-scoring algorithm, where the Fisher information matrix is obtained using (13), (14) and (15),

$$\mathbf{I} = \begin{pmatrix} \mathbf{I}_{\tau_1 \tau_1} & \mathbf{I}_{\tau_1 \tau_2} & \mathbf{I}_{\tau_1 \boldsymbol{\theta}} \\ \mathbf{I}_{\tau_2 \tau_1} & \mathbf{I}_{\tau_2 \tau_2} & \mathbf{I}_{\tau_2 \boldsymbol{\theta}} \\ \mathbf{I}_{\boldsymbol{\theta} \tau_1} & \mathbf{I}_{\boldsymbol{\theta} \tau_2} & \mathbf{I}_{\boldsymbol{\theta} \boldsymbol{\theta}} \end{pmatrix} = \begin{pmatrix} \mathbf{I}_{\tau_1 \tau_1}^T & \mathbf{I}_{\tau_1 \tau_2} & \mathbf{I}_{\tau_1 \boldsymbol{\theta}} \\ \mathbf{I}_{\tau_2 \tau_1}^T & \mathbf{I}_{\tau_2 \tau_2} & \mathbf{I}_{\tau_2 \boldsymbol{\theta}} \\ \mathbf{I}_{\boldsymbol{\theta} \tau_1}^T & \mathbf{I}_{\boldsymbol{\theta} \tau_2} & \mathbf{I}_{\boldsymbol{\theta} \boldsymbol{\theta}} \end{pmatrix}$$

where

$$\mathbf{I}_{\tau_1 \tau_1} = \frac{1}{2} \text{Tr}(\mathbf{P}_* \mathbf{B}_{1*} \mathbf{B}_{1*}^T \mathbf{P}_* \mathbf{B}_{1*} \mathbf{B}_{1*}^T), \quad \mathbf{I}_{\tau_2 \tau_2} = \frac{1}{2} \text{Tr}(\mathbf{P}_* \mathbf{B}_{2*} \mathbf{B}_{2*}^T \mathbf{P}_* \mathbf{B}_{2*} \mathbf{B}_{2*}^T),$$

$$\mathbf{I}_{\tau_1 \boldsymbol{\theta}_j} = \frac{1}{2} \text{Tr} \left(\mathbf{P}_* \mathbf{B}_{1*} \mathbf{B}_{1*}^T \mathbf{P}_* \frac{\partial \mathbf{V}}{\partial \theta_j} \right), \quad \mathbf{I}_{\tau_2 \boldsymbol{\theta}_j} = \frac{1}{2} \text{Tr} \left(\mathbf{P}_* \mathbf{B}_{2*} \mathbf{B}_{2*}^T \mathbf{P}_* \frac{\partial \mathbf{V}}{\partial \theta_j} \right),$$

and

$$\mathbf{I}_{\tau_1 \tau_2} = \frac{1}{2} \text{Tr}(\mathbf{P}_* \mathbf{B}_{1*} \mathbf{B}_{1*}^T \mathbf{P}_* \mathbf{B}_{2*} \mathbf{B}_{2*}^T), \quad \mathbf{I}_{\boldsymbol{\theta}_j \boldsymbol{\theta}_k} = \frac{1}{2} \text{Tr} \left(\mathbf{P}_* \frac{\partial \mathbf{V}}{\partial \theta_j} \mathbf{P}_* \frac{\partial \mathbf{V}}{\partial \theta_k} \right).$$

4 Simulation Study

We conduct a toy simulation study to evaluate the performance of the estimation procedure of the model regression parameters and nonparametric function using the REML estimates for the smoothing parameters and the variance parameters. Bivariate cyclic longitudinal data are generated according to the following model:

$$\begin{aligned} Y_{1ij} &= \text{age}_i^T \beta_1 + f_1(t_{ij}) + b_{1i} + U_{1i}(t_{ij}) + \epsilon_{1ij} \\ Y_{2ij} &= \text{age}_i^T \beta_2 + f_2(t_{ij}) + b_{2i} + U_{2i}(t_{ij}) + \epsilon_{2ij} \\ &\quad i = 1, \dots, 30; \quad j = 1, \dots, 28; \quad t_{ij} \in \{1, \dots, 28\} \end{aligned}$$

where b_{1i} and b_{2i} are independent but correlated random intercepts following a bivariate normal distribution:

$$\begin{pmatrix} b_{1i} \\ b_{2i} \end{pmatrix} \sim \mathbf{N}_2 \left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} \phi_1 & \phi_2 \\ \phi_2 & \phi_3 \end{pmatrix} \right);$$

U_{1i} and U_{2i} are simulated from mean 0 bivariate NOU fields modeling serial correlation, with variance function $\text{var}(U_{1i}(t)) = \exp\{a_{10} + a_{11}t + a_{12}t^2\}$, $\text{var}(U_{2i}(t)) = \exp\{a_{20} + a_{21}t + a_{22}t^2\}$ and $\text{corr}(U_{1i}(t), U_{1i}(s)) = \rho_1^{|s-t|}$ $\text{corr}(U_{2i}(t), U_{2i}(s)) = \rho_2^{|s-t|}$, i.e. the covariance function for the bivariate NOU field is

$$\mathbf{C}_i(s, t) = \begin{pmatrix} \rho_1^{|s-t|} \exp\{a_{10} + a_{11}t + a_{12}t^2\} & 0 \\ 0 & \rho_2^{|s-t|} \exp\{a_{20} + a_{21}t + a_{22}t^2\} \end{pmatrix};$$

lastly, ϵ_{1ij} and ϵ_{2ij} are simulated from a mean 0 bivariate normal distribution

$$\begin{pmatrix} \epsilon_{1i} \\ \epsilon_{2i} \end{pmatrix} \sim \mathbf{N}_2 \left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} \sigma_1^2 & 0 \\ 0 & \sigma_2^2 \end{pmatrix} \right).$$

Further, the nonparametric functions are generated from

$$f_1(t) = 5 \sin \left(\frac{2\pi}{28} \right) t, \quad f_2(t) = 3 \cos \left(\frac{2\pi}{28} \right) t$$

with periods to be 28 days for both responses.

Table 1: Simulation results for estimates of model parameters based on 100 simulation replicates

| Model parameters | True Value | Parameter estimate | Bias | SE |
|------------------|------------|--------------------|---------|--------|
| β_1 | 1.0000 | 0.9979 | 0.0021 | 0.0017 |
| β_2 | 0.7500 | 0.7502 | 0.0003 | 0.0018 |
| τ_1 | 1.0000 | 0.5652 | 0.4348 | 0.0095 |
| τ_2 | 1.0000 | 0.5939 | 0.4061 | 0.0100 |
| ϕ_1 | 1.0000 | 0.9949 | 0.0051 | 0.0107 |
| ϕ_2 | -0.5000 | -0.5092 | -0.0184 | 0.0085 |
| ϕ_3 | 1.0000 | 0.9939 | 0.0061 | 0.0114 |
| σ_1^2 | 1.0000 | 1.0008 | 0.0008 | 0.0028 |
| σ_2^2 | 1.0000 | 0.9991 | 0.0009 | 0.0029 |
| ρ_1 | 0.9163 | 0.0831 | 0.9093 | 0.0027 |
| a_{10} | -5.0000 | -5.0360 | -0.0072 | 0.0808 |
| a_{11} | 1.5000 | 1.5132 | 0.0088 | 0.0216 |
| a_{12} | -0.1000 | -0.1012 | -0.0120 | 0.0014 |
| ρ_2 | 0.9163 | 0.0854 | 0.9068 | 0.0027 |
| a_{20} | -5.0000 | -4.9785 | -0.0043 | 0.0771 |
| a_{21} | 1.5000 | 1.4933 | 0.0045 | 0.0207 |
| a_{22} | -0.1000 | -0.0998 | -0.0020 | 0.0014 |

Table 1 recorded the simulation results for estimates of model parameters based on 500 simulation replicates and 30 subjects. The Bias is defined as the bias of the parameter estimated divided by its true value. The parameter estimates of the regression coefficients β_1 and β_2 , and the variance estimates of the random intercepts and measurement errors are nearly unbiased, whereas the estimates of the smoothing parameters and the NOU variance parameters are slightly biased.

The biases for the nonparametric functions \hat{f}_1 and \hat{f}_2 are minimal for \hat{f}_2 which centers around 0, whereas for \hat{f}_1 the bias stands a little above 0, see Figure 4. Figure 4 shows that model standard errors of estimates of \hat{f}_1 and \hat{f}_2 agree quite well with the empirical standard errors.

Figure 4 shows the estimated pointwise 95% coverage probabilities of the true nonparametric functions f_1 and f_2 . The means for the estimated coverage probabilities are 96% and 93% for \hat{f}_1 and \hat{f}_2 . Our simulation results are largely consistent with those of Zhang et al. (1998) for univariate independent data.

5 SWAN data analysis

The model we proposed was motivated by a dataset from the Study of Women's Health Across the Nation (SWAN), a cohort study of 3,302 middle-aged women in the United States, collected from seven sites. The women enrolled in the study come from different ethnic backgrounds. Daily urine samples were collected from 403 employed women aged 20 to 44 years who completed a median of five consecutive menstrual cycles of collection each [2]. Of these, 338 women collected daily urine samples for at least one complete menstrual cycle, had fewer than three days of missing data in any five-day rolling window, did not have a conception in the analyzed cycles, and had complete covariate information. One menstrual cycle was randomly selected

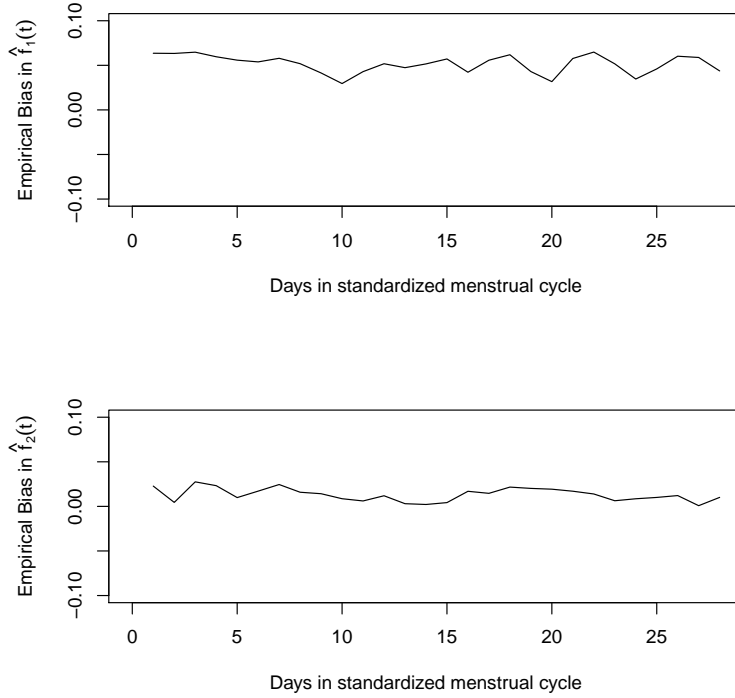


Figure 1: Empirical Bias in estimated nonparametric functions \hat{f}_1 and \hat{f}_2 based on 500 simulation replications.

from each of the 338 women. Risk factor data were obtained by in-person interview at baseline. The details of the study design and assay methods have been described in detail previously [2] [3].

We are interested in modelling the mean curve for women’s daily urinary estrogen (EIC) and progesterone (PdG) metabolite profiles, and their relationships to demographic and lifestyle factors over a 28-day reference menstrual cycle. Also, we are interested in their potential interactions between these two hormones. Thus, we will model jointly for these two responses. For demonstration purposes, we randomly select 50 study participants from the study, with a total of 2714 observations for both responses. Each woman contributes from 16 to 41 observations over a menstrual cycle, resulting an average of 27 observations per woman. In order for the results to be biologically meaningful, the menstrual cycle length for each woman has been standardized to a reference of 28 days, based on the assumption that the change of hormone level for each woman depends on the time of the menstrual cycle relative to the cycle length. The standardization generates 56 distinct time points. To make the normality assumption more appropriate, the log transformation was used for both responses.

Figure (5) plots the log-transformed progesterone and estrogen levels during a standardized menstrual cycle. Figure (5) plots their empirical sample variances calculated at each distinct time points.

Denote $\{(Y_{1ij}, Y_{2ij})\}$ the j^{th} log-transformed progesterone and estrogen values measured at standardized day t_{ij} since menstruation for the i^{th} woman, we consider the following bivariate semiparametric stochastic mixed model:

$$\begin{aligned} Y_{1ij} &= \text{age}_i^T \beta_{11} + \text{BMI}_i^T \beta_{12} + f_1(t_{ij}) + b_{1i} + U_i(t_j) + \epsilon_{1ij} \\ Y_{2ij} &= \text{age}_i^T \beta_{21} + \text{BMI}_i^T \beta_{22} + f_2(t_{ij}) + b_{2i} + U_i(t_j) + \epsilon_{2ij} \\ i &= 1, \dots, 50; j = 1, \dots, n_i; t_{ij} \in \{0.5, 1.0, \dots, 28\} \end{aligned}$$

where b_{1i} and b_{2i} are the random intercepts that are correlated between the two hormone response but inde-

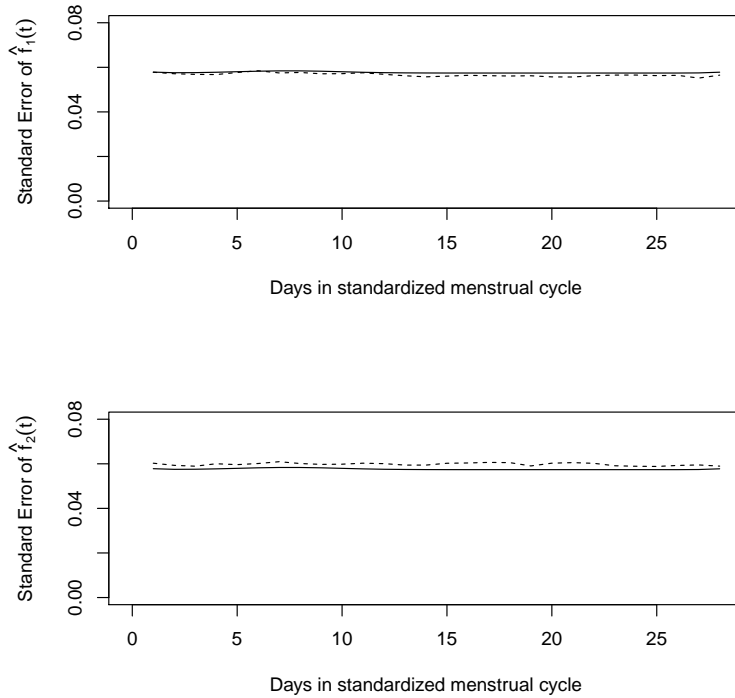


Figure 2: Pointwise standard errors of the estimated nonparametric functions \hat{f}_1 and \hat{f}_2 based on 500 simulation replications.

pendent between subjects following a bivariate normal distribution with mean zero and variance $N_2 \left(\mathbf{0}, \begin{pmatrix} \phi_1 & \phi_3 \\ \phi_3 & \phi_2 \end{pmatrix} \right)$; the U_i are mean 0 bivariate Gaussian field modeling serial correlation, and the ϵ_{1ij} and ϵ_{2ij} are independent mean zero measurement errors following bivariate normal with variance $\begin{pmatrix} \sigma_1^2 & 0 \\ 0 & \sigma_2^2 \end{pmatrix}$.

| Model parameters | Parameter estimate | SE | CI |
|------------------|--------------------|--------|-------------------|
| β_{11} | -0.0127 | 0.0243 | (-0.0603, 0.0349) |
| β_{12} | -0.1181 | 0.1973 | (-0.5048, 0.2686) |
| β_{21} | -0.0191 | 0.0224 | (-0.0630, 0.0248) |
| β_{22} | 0.0857 | 0.1814 | (-0.2698, 0.4412) |

6 Conclusion and Discussions

In conclusion, we propose and build a model for bivariate cyclic longitudinal data. The model is proposed in the likelihood framework and the regression parameters and nonparametric functions are estimated by maximizing penalized likelihood function. The smoothing parameter and variance components are numerically estimated using the Fisher-scoring algorithm based on restricted maximum likelihood. Modelling the time effect nonparametrically gives more flexibility. The Gaussian field allows for additional flexibility in within-subject correlation structure than that generated in random effects.

The model we proposed can also be readily extended to multivariate cyclic longitudinal data. One thing to note is that dimensionality can pose as a challenge when extending to multivariate cyclic longitudinal data. In the bivariate studies, we employed both C++ and parallel computing in the simulation study. Despite

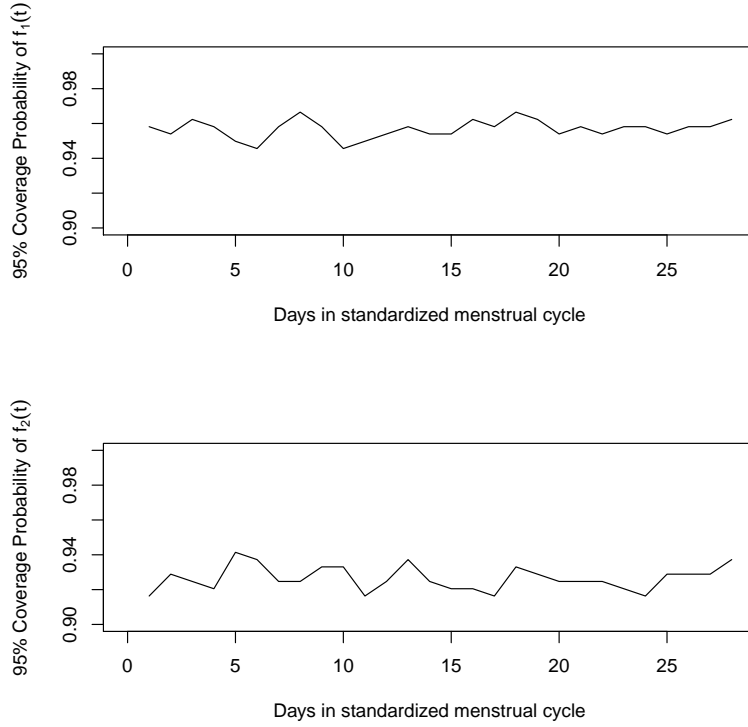


Figure 3: Estimated pointwise 95% coverage probabilities of the true nonparametric functions f_1 and f_2 based on 500 simulation replications.

the effort, the computation time is still nontrivial. We can therefore infer that extension to multivariate model will give rise to significant computation challenge.

For future work, we would like to extend an existing univariate Bayesian model for cyclic longitudinal data to build a model for the bivariate cyclic longitudinal data under the Bayesian’s framework, and apply both frequentist model and Bayesian model to the SWAN data and compare results. Also, we will explore the sensitivity and robustness to the model assumption.

References

- [1] Babette A. Brumback and John A. Rice. Smoothing spline models for the analysis of nested and crossed samples of curves. *Journal of the American Statistical Association*, 93(443):961–976, 1998.
- [2] Ellen B Gold, Brenda Eskenazi, Bill L Lasley, Steven J Samuels, Marianne O’Neill Rasor, James W Overstreet, and Marc B Schenker. Epidemiologic methods for prospective assessment of menstrual cycle and reproductive characteristics in female semiconductor workers. *American journal of industrial medicine*, 28(6):783–797, 1995.
- [3] Ellen B Gold, Brenda Eskenazi, S Katharine Hammond, Bill L LaSley, Steven J Samuels, Marianne O’Neill Rasor, Cynthia J Hines, James W Overstreet, and Marc B Schenker. Prospectively assessed menstrual cycle characteristics in female wafer-fabrication and nonfabrication semiconductor employees. *American journal of industrial medicine*, 28(6):799–815, 1995.
- [4] P. J. Green. Penalized likelihood for general semi-parametric regression models. *International Statistical Review*, 55:245–260, 1987.

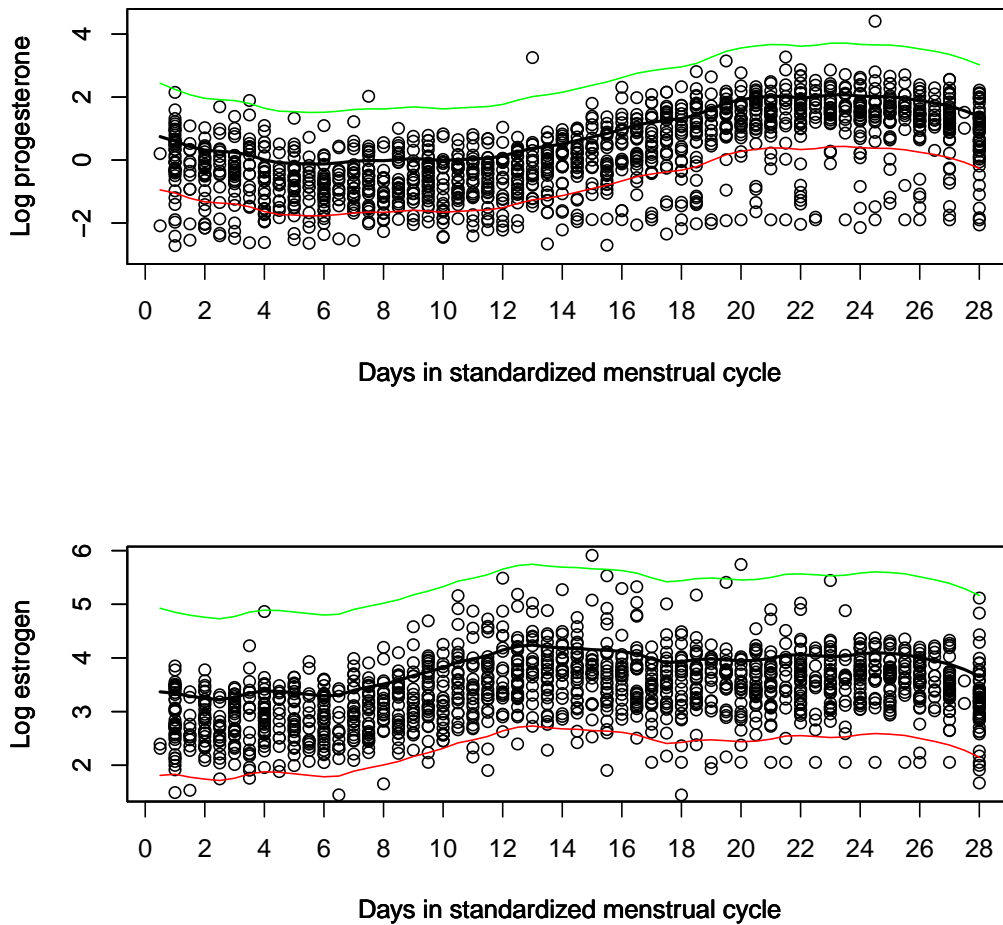


Figure 4: Plots of log progesterone and log estrogen levels against days in a standardized menstrual cycle.

- [5] P. J. Green and Bernard W. Silverman. *Nonparametric Regression and Generalized Linear Models: A roughness penalty approach*. Chapman and Hall/CRC, London, UK, 1994.
- [6] Leonid B. Koralov and Yakov G. Sinai. *Theory of Probability and Random Processes*. Springer, New York, NY, second edition, 2007.
- [7] Ziyue Liu, Anne R. Cappola, Leslie J. Crofford, and Wensheng Guo. Modeling bivariate longitudinal hormone profiles by hierarchical state space models. *Journal of the American Statistical Association*, 109(505):108–118, 2014.
- [8] P. M. Meyer, S. L. Zeger, Sioban D. Harlow, M. Sowers, S. Crawford, J. L. Luborsky, I. Janssen, D. S. McConnell, J. F. Randolph, and G. Weiss. Characterizing daily urinary hormone profiles for women at midlife using functional data analysis. *American Journal of Epidemiology*, 165(8):936–945, 2007.
- [9] J. P. Sy, J. M. G. Taylor, and W. G. Cumberland. A stochastic model for the analysis of bivariate longitudinal aids data. *Biometrics*, 53(2):542–555, 1997.

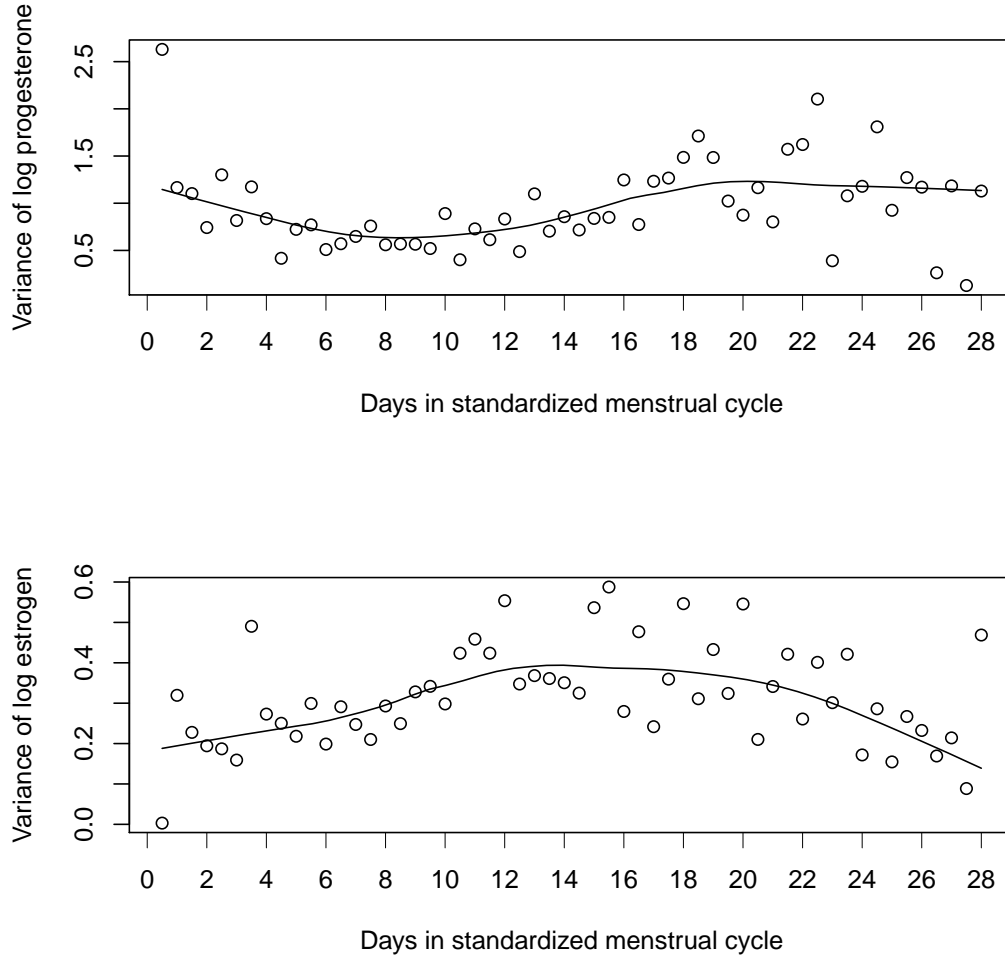


Figure 5: Plots of empirical sample variance of log progesterone and log estrogen levels at each distinct time points in a standardized menstrual cycle.

- [10] J. M. G. Taylor, W. G. Cumberland, and J. P. Sy. A stochastic model for analysis of longitudinal aids data. *Journal of the American Statistical Association*, 89:727–736, 1994.
- [11] Arunas P. Verbyla, Brian R. Cullis, Michael G. Kenward, and Sue J. Welham. The analysis of designed experiments and longitudinal data by using smoothing splines. *Journal of the Royal Statistical Society. Series C (Applied Statistics)*, 48(3):269–311, 1999.
- [12] Yuedong Wang, Wensheng Guo, and Morton B Brown. Spline smoothing for bivariate data with applications to association between hormones. *Statistica Sinica*, pages 377–397, 2000.
- [13] Sue J. Welham, Brian R. Cullis, Michael G. Kenward, and Robin Thompson. The analysis of longitudinal data using mixed model l-splines. *Biometrics*, 62(2):392–401, 2006.
- [14] Simon Wood. *Generalized Additive Models: An Introduction with R*. Chapman and Hall/CRC, Boca Raton, Florida, 2006.

- [15] D. Zhang, X. Lin, and M. Sowers. Semiparametric regression for periodic longitudinal hormone data from multiple menstrual cycles. *Biometrics*, 56(1):31–39, 2000.
- [16] Daowen Zhang, Xihong Lin, Jonathan Raz, and MaryFran Sowers. Semiparametric stochastic mixed models for longitudinal data. *Journal of the American Statistical Association*, 93(442):710–719, 1998.