

# Fixed and random effects selection in nonparametric additive mixed models

Randy C. S. Lai

*Department of Statistics, University of California at Davis  
4118 Mathematical Sciences Building, One Shields Avenue  
Davis, CA 95616, USA  
e-mail: [rcslai@ucdavis.edu](mailto:rcslai@ucdavis.edu)*

Hsin-Cheng Huang\*

*Institute of Statistical Science  
Academia Sinica, Taipei 115, Taiwan  
e-mail: [hchuang@stat.sinica.edu.tw](mailto:hchuang@stat.sinica.edu.tw)*

and

Thomas C. M. Lee<sup>†</sup>

*Department of Statistics, University of California at Davis  
4118 Mathematical Sciences Building, One Shields Avenue  
Davis, CA 95616, USA  
e-mail: [tcmLee@ucdavis.edu](mailto:tcmLee@ucdavis.edu)*

**Abstract:** This paper considers the problem of model selection in a nonparametric additive mixed modeling framework. The fixed effects are modeled nonparametrically using truncated series expansions with B-spline basis. Estimation and selection of such nonparametric fixed effects are simultaneously achieved by using the adaptive group lasso methodology, while the random effects are selected by a traditional backward selection mechanism. To facilitate the automatic selection of model dimension, computable expressions for the degrees of freedom for both the fixed and random effects components are derived, and the Bayesian Information criterion (BIC) is used to select the final model choice. Theoretically it is shown that this BIC model selection method is consistent, while computationally a practical algorithm is developed for solving the optimization problem involved. Simulation results show that the proposed methodology is often capable of selecting the correct significant fixed and random effects components, especially when the sample size and/or signal to noise ratio are not too small. The new method is also applied to two real data sets.

**AMS 2000 subject classifications:** Primary 62G08.

**Keywords and phrases:** Adaptive group lasso, additive mixed model, Bayesian information criterion, consistency.

Received March 2011.

---

\*Supported in part by National Science Council Taiwan under grants 97-2118-M001-001-MY3 and 100-2628-M-001-004-MY3.

<sup>†</sup>Supported in part by NSF under Grant DMS 1007520.

## 1. Introduction

This paper considers the problem of joint selection of fixed and random effects in nonparametric additive mixed models. Suppose for the  $i$ -th subject we observe response  $Y_i$  and covariates  $X_{i1}, \dots, X_{iP}$  and  $Z_{i1}, \dots, Z_{iQ}$ . It is assumed that these measurements are related by the following mixed model:

$$Y_i = \mu + \sum_{p=1}^P f_p(X_{ip}) + \sum_{q=1}^Q u_q Z_{iq} + \epsilon_i, \quad i = 1, \dots, n, \quad (1)$$

where  $\mu$  is the mean,  $f_p$ 's are unknown smooth continuous functions that form the fixed component of the model,  $\mathbf{u} = (u_1, \dots, u_Q)'$  are random effects, and  $\boldsymbol{\epsilon} = (\epsilon_1, \dots, \epsilon_n)'$  are additive errors. For simplicity, write  $\mathbf{Y} = (Y_1, \dots, Y_n)'$ , and define  $\mathbf{Z}$  as the matrix with  $Z_{iq}$  as its  $iq$ -th element, for  $i = 1, \dots, n$  and  $q = 1, \dots, Q$ . It is further assumed that  $\mathbf{u}$  and  $\boldsymbol{\epsilon}$  follow the normal distribution and satisfy

$$E \begin{pmatrix} \mathbf{u} \\ \boldsymbol{\epsilon} \end{pmatrix} = \mathbf{0}, \quad \text{Cov} \begin{pmatrix} \mathbf{u} \\ \boldsymbol{\epsilon} \end{pmatrix} = \sigma^2 \begin{pmatrix} \mathbf{G}(\boldsymbol{\theta}) & \mathbf{0} \\ \mathbf{0} & \mathbf{I} \end{pmatrix}$$

and

$$\mathbf{V} = \text{Cov}(\mathbf{Z}\mathbf{u} + \boldsymbol{\epsilon}) = \sigma^2(\mathbf{Z}\mathbf{G}\mathbf{Z}' + \mathbf{I}),$$

where  $\sigma^2 > 0$  and  $\mathbf{G}(\boldsymbol{\theta})$  is a positive definite matrix with elements known upto the parameter vector  $\boldsymbol{\theta}$ ; in the sequel we simply write  $\mathbf{G} = \mathbf{G}(\boldsymbol{\theta})$ . For identifiability purposes, we assume  $E[f_p(X_{ip})] = 0$  for all  $i$  and  $p$ . We consider the case that the  $X_{ip}$ 's are random and sampled from some continuous density function.

The concept of fixed and random effects has been successfully applied to repeated measurement data to study the variations both within and between replicates of subjects; e.g., see Laird and Ware (1982); Zeger and Liang (1986); Liang and Zeger (1986); Diggle *et al.* (2002). These authors mainly consider the case when the fixed effects are linear, while our approach can be seen as an extension in which nonparametric fixed effects are allowed. The nonparametric additive mixed model for repeated measurements data can be written as:

$$Y_{ik} = \mu + \sum_{p=1}^P f_p(X_{ikp}) + \sum_{q=1}^Q u_{iq} Z_{ikq} + \epsilon_{ik}, \quad i = 1, \dots, N, \quad k = 1, \dots, n_i, \quad (2)$$

where  $Y_{ik}$  is the  $k$ th measurement of the  $i$ th subject,  $X_{ikp}$  and  $Z_{ikq}$  are covariates associated with  $Y_{ik}$ ,  $N$  is the total number of subjects and  $n_i$  is the number of measurements taken from the  $i$ th subject. The random effects  $u_{iq}$ 's model the variation among different subjects while the fixed effects  $f_p$ 's model the variation between the replicates. Equation (2) can be transformed to (1) by suppressing the replicate index  $k$ , and the total number of observations is  $\sum_{i=1}^N n_i = n$ . Similar formulations of modeling have been discussed for examples in Zhang

*et al.* (1998); Lin and Zhang (1999); Wand (2003); Fahrmeir and Lang (2001). The first three pieces of work consider  $L_2$  minimization criteria while the last considers the problem from a Bayesian perspective.

The problem of fixed and random effects selection has been well studied for linear mixed models. Chen and Dunson (2003) develop a Bayesian approach for selecting random effects in linear mixed models, while Kinney and Dunson (2007) use Bayesian method for selecting both fixed and random effects in linear and logistic mixed models. Bondell *et al.* (2010) also study the problem of fixed and random effects selection in linear mixed models, but with a penalized likelihood approach. Lastly, Pu and Niu (2006) extended the generalized information criterion for choosing fixed and random effects in linear mixed models.

Given (1) (or (2)), the goal of this article is to identify which of the functions  $f_p$ 's and which of the random effects  $u_q$ 's are statistically significant. Our methodology begins with modeling the  $f_p$ 's nonparametrically by B-spline basis expansions. As to be demonstrated in Section 2 below, the use of B-spline bases can transform (1) into a compact matrix representation. For parameter estimation, we develop an iterative algorithm that combines the adaptive group lasso and Newton-Raphson methods. This algorithm is presented in Section 3. Section 4 derives a Bayesian information criterion to perform simultaneous variable selection for both fixed and random effects. Theoretical properties of our methodology are studied in Section 5. We illustrate the practical performance of our approach via both simulation experiments and practical data analysis; see Sections 6 and 7. Technical details are deferred to the appendix.

## 2. B-spline modeling of nonparametric fixed effects

Loosely, a spline is a piecewise polynomial that is smoothly connected at its knots. Denote  $\mathcal{S}(d_0, \mathbf{t})$  as the collection of all spline functions of order  $d_0$  with knots  $\mathbf{t} = (t_0, \dots, t_{\tau+1})'$ , where  $t_0 \leq t_1 \leq \dots \leq t_{\tau+1}$ . There exists  $m = d_0 + \tau$  B-spline basis functions  $\{\phi_j(\cdot)\}_{j=1}^m$  for  $\mathcal{S}(d_0, \mathbf{t})$ . That is, all elements of  $\mathcal{S}(d_0, \mathbf{t})$  can be expressed as linear combinations of  $\{\phi_j(\cdot)\}_{j=1}^m$ . These B-spline basis functions are normalized in the sense that  $\sum_{j=1}^m \phi_j(x) = 1$  for all  $x$ . They play an important role in non-parametric additive modeling; e.g., see Stone (1985, 1986). For exact expressions of  $\phi_j$ 's, see de Boor (2001).

Denote  $\mathbf{t}_p$  as the knot vector for  $f_p$ , and let  $\{\phi_{pj}(\cdot)\}_{j=1}^{m_p}$  be the corresponding basis functions of  $\mathcal{S}(d_0, \mathbf{t}_p)$ . We model  $f_p$  by  $\tilde{f}_p$ :

$$\tilde{f}_p = \sum_{k=1}^{m_p} \beta_{pk} \phi_{pk}.$$

To make the model identifiable, we assume  $\sum_{i=1}^n \tilde{f}_p(X_{ip}) = 0$ , which is a sample analogue of  $E[f_p(X_{ip})] = 0$ . This can be achieved by centering the basis functions, or equivalently, by assuming that

$$\sum_{k=1}^{m_p} \sum_{i=1}^n \phi_{pk}(X_{ip}) \beta_{pk} = 0. \quad (3)$$

Hence, the number of effective basis functions is  $m_p$ . For simplicity, we set  $m_p = m$  for all  $p$ ; the effect of this is minimal as long as  $m$  is large enough.

For  $p = 1, \dots, P$ , let  $\boldsymbol{\beta}_p = (\beta_{p1}, \dots, \beta_{pm})'$  and  $\boldsymbol{\beta} = (\mu, \boldsymbol{\beta}'_1, \dots, \boldsymbol{\beta}'_P)'$ . Also, write  $\mathbf{X}_p$  as the matrix with its  $ij$ -th element as  $\phi_{pj}(X_{ip})$  for  $j = 1, \dots, m$  and  $i = 1, \dots, n$ , and define  $\mathbf{X} = [\mathbf{1}, \mathbf{X}_1, \dots, \mathbf{X}_P]$ . The matrix representation of our spline model for (1) is

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u} + \boldsymbol{\epsilon} = \mu\mathbf{1} + \sum_{p=1}^P \mathbf{X}_p\boldsymbol{\beta}_p + \mathbf{Z}\mathbf{u} + \boldsymbol{\epsilon}, \quad (4)$$

which is a common formulation of mixed models. However, in (4) those columns of  $\mathbf{X}$  that correspond to the basis functions of the same  $f_p$  are naturally grouped together, in the sense that this whole group of basis functions should be kept or killed together during the fixed effects selection process. This concept of group variables has been considered in Lin and Zhang (2006); Yuan and Lin (2006); Wang *et al.* (2007b). We note that (4) can also be expressed as  $\mathbf{Y} \sim N(\mathbf{X}\boldsymbol{\beta}, \mathbf{V})$ , with the log-likelihood function given as

$$l(\boldsymbol{\beta}, \mathbf{V}) = -\frac{1}{2} \log(|\mathbf{V}|) - \frac{1}{2} (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})' \mathbf{V}^{-1} (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}).$$

### 3. Parameter estimation

There are two types of parameters in (4): the fixed component parameters  $\boldsymbol{\beta} = (\mu, \boldsymbol{\beta}'_1, \dots, \boldsymbol{\beta}'_P)$  and the random component parameters  $\{\mathbf{G}, \sigma^2\}$ . This section develops a method for estimating these parameters when the complexity of the model is pre-specified. Automatic choice of model complexity will be discussed in the next section. The main idea of our estimation method is to iterate the following two steps until convergence:

1. given a current set of estimates for  $\{\mathbf{G}, \sigma^2\}$ , estimate  $\boldsymbol{\beta}$  via the adaptive group lasso methodology,
2. given a current estimate for  $\boldsymbol{\beta}$ , obtain estimates for  $\{\mathbf{G}, \sigma^2\}$  using the Newton-Raphson method.

#### 3.1. Fixed component estimation using adaptive group lasso

The lasso was introduced by Tibshirani (1996) as a methodology for simultaneous variable selection and shrinkage estimation. Later Zou (2006) developed a variant termed the adaptive lasso which was shown to possess superior theoretical properties. The use of lasso for grouped variables, the so-called group lasso, has been discussed for examples in Yuan and Lin (2006); Meier *et al.* (2008); Wei and Huang (2008); Ravikumar *et al.* (2009). In particular, Huang *et al.* (2010); Meier *et al.* (2009) consider the estimation of additive models using group lasso. Here we develop an adaptive group lasso procedure for estimating the fixed component parameters.

Suppose for now estimates for  $\{\mathbf{G}, \sigma^2\}$  and hence  $\mathbf{V}$  are available. Denote the estimate as  $\hat{\mathbf{V}}$ , and write the likelihood as  $l(\boldsymbol{\beta}) = l(\boldsymbol{\beta}, \hat{\mathbf{V}})$ . Firstly the estimate for  $\mu$  is  $\hat{\mu} = \bar{Y}$ , the average of all  $Y_i$ 's. Then our adaptive group lasso estimates for the remaining parameters in  $\boldsymbol{\beta} = (\mu, \boldsymbol{\beta}'_1, \dots, \boldsymbol{\beta}'_P)$  are defined as the minimizer of

$$S(\boldsymbol{\beta}) = -l(\boldsymbol{\beta}) + \lambda \sum_{p=1}^P w_p \|\boldsymbol{\beta}_p\| \quad (5)$$

with the weights  $w_p$  given by

$$w_p = \begin{cases} (\|\tilde{\boldsymbol{\beta}}_p\|^{-1} - \nu)_+, & \text{if } \|\tilde{\boldsymbol{\beta}}_p\| > 0, \\ \infty, & \text{otherwise,} \end{cases} \quad p = 1, \dots, P, \quad (6)$$

where  $\|\cdot\|$  is the usual Euclidean norm,  $\lambda > 0$  and  $\nu > 0$  are tuning parameters, and  $\tilde{\boldsymbol{\beta}}_p$  is an initial estimate for  $\boldsymbol{\beta}_p$ . Both  $\lambda$  and  $\nu$  are pre-specified before the minimization of (5) is carried out. When  $\lambda$  is zero, the corresponding estimates will be reduced to the generalized least square estimate. As  $\lambda$  becomes larger, more  $\boldsymbol{\beta}_p$ 's are shrunk to zero. The presence of  $\nu$  is to ensure the consistency properties of our estimation method (see Section 5). Methods for choosing  $\lambda$  and  $\nu$  will be provided in Section 4 below.

Due to the non-differentiable property of the penalty term, minimization of (5) is difficult. Different algorithms have been developed for minimizing such group lasso type likelihoods; e.g., Meier *et al.* (2008); Yuan and Lin (2006). For (5), as the B-spline basis functions are not orthogonal and the variance component  $\mathbf{V}$  is not diagonal, the shooting algorithm of Yuan and Lin (2006) is not applicable, but the blockwise co-ordinate gradient descent (BCGD) algorithm of Meier *et al.* (2008) can be adopted. Below is a brief description of the BCGD algorithm customized to our settings.

Note that the log-likelihood  $l(\boldsymbol{\beta})$  is not a quadratic function and the key idea of BCGD is to combine a quadratic approximation of  $l(\boldsymbol{\beta})$  with a line search. First, for any  $\mathbf{d} \in \mathbb{R}^{mP+1}$ , define the following approximation  $M(\mathbf{d})$  to  $S(\boldsymbol{\beta} + \mathbf{d})$ :

$$M(\mathbf{d}) = - \left\{ l(\boldsymbol{\beta}) + \mathbf{d}' \frac{\partial l(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}} + \frac{1}{2} \mathbf{d}' \mathbf{H} \mathbf{d} \right\} + \lambda \sum_{p=1}^P w_p \|\boldsymbol{\beta}_p + \mathbf{d}_p\| \approx S(\boldsymbol{\beta} + \mathbf{d}),$$

where  $\mathbf{H}$  is a diagonal matrix approximating the Hessian of  $S(\boldsymbol{\beta} + \mathbf{d})$ . Denote the diagonal elements of  $\mathbf{H}$  as  $(h_0, h_1 \mathbf{e}_m, h_2 \mathbf{e}_m, \dots, h_p \mathbf{e}_m, \dots, h_P \mathbf{e}_m)$ , where  $\mathbf{e}_m$  is a row vector of  $m$  ones. For  $p > 0$ , a possible choice of  $h_p$  is

$$h_p = \min \left[ \text{diag} \left\{ \frac{\partial^2 l(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}_p \partial \boldsymbol{\beta}'_p} \right\} \right] = - \max[\text{diag}(\mathbf{X}'_p \mathbf{V}^{-1} \mathbf{X}_p)].$$

For  $h_0$ , it vanishes after a differentiation operation and hence plays no role in the actual minimization algorithm. For simplicity, we set  $h_0 = 0$ .

Next we consider minimizing  $M(\mathbf{d})$  one “block at a time”. More precisely, for each  $p > 0$ , we minimize  $M(\mathbf{d})$  with respect to (w.r.t.)

$$\mathbf{d} = (0, 0\mathbf{e}_m, \dots, 0\mathbf{e}_m, \mathbf{d}_p, 0\mathbf{e}_m, \dots, 0\mathbf{e}_m),$$

where  $\mathbf{d}_p \in \mathbb{R}^m$ . Direct algebra shows that

$$\mathbf{d}_p = \begin{cases} -\boldsymbol{\beta}_p, & \text{if } \|\mathbf{S}_p\| < \lambda w_p \\ -\frac{1}{h_p} \left\{ \frac{\partial l(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}_p} - \lambda w_p \frac{\mathbf{S}_p}{\|\mathbf{S}_p\|} \right\}, & \text{otherwise} \end{cases}, \quad (7)$$

where  $\mathbf{S}_p = \frac{\partial l(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}_p} - h_p \boldsymbol{\beta}_p$ . Therefore, given  $\boldsymbol{\beta}^{(t)}$ ,  $\mathbf{d}_p^{(t)}$  is the best direction of search w.r.t. the  $p$ th component for minimizing  $S(\boldsymbol{\beta}^{(t)})$ .

Lastly, we update the minimizer by  $\boldsymbol{\beta}^{(t+1)} + a^{(t)} \mathbf{d}_p^{(t)}$ , where  $a^{(t)}$  is the step length of the search. This step length can be chosen by the Armijo rule:  $a^{(t)}$  is chosen as the largest value that satisfies

$$S(\boldsymbol{\beta}^{(t+1)}) \leq S(\boldsymbol{\beta}^{(t)}) + c_1 a_p^{(t)} \Delta^{(t)}, \quad (8)$$

where  $\Delta^{(t)} = -\mathbf{d}_p^{(t)'} \frac{\partial l(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}} + \lambda \{ \|\boldsymbol{\beta}_p^{(t)} + \mathbf{d}_p^{(t)}\| - \|\boldsymbol{\beta}^{(t)}\| \}$  and  $c_1 > 0$ . The value  $\Delta^{(t)}$  is a linear approximation of the improvement in  $S(\cdot)$  and we choose  $c_1 = 0.1$  as suggested by Meier *et al.* (2008). Selecting the step length  $a^{(t)}$  in this fashion ensures a sufficient decrease in the penalized likelihood.

In summary, given a current estimate  $\hat{\boldsymbol{\beta}}^{(t)}$  of  $\boldsymbol{\beta}$ , the BCGD algorithm computes the next iterative estimate  $\hat{\boldsymbol{\beta}}^{(t+1)}$  with the following steps: for  $p = 1, \dots, P$ ,

1. evaluates  $\mathbf{d}_p^{(t)}$  with (7),
2. finds the largest  $a^{(t)}$  that satisfies (8), and
3. set  $\hat{\boldsymbol{\beta}}_p^{(t+1)} = \hat{\boldsymbol{\beta}}_p^{(t)} + a_p^{(t)} \mathbf{d}_p^{(t)}$ .

This BCGD algorithm converged rapidly in all our numerical work.

### 3.2. Random component estimation using Newton Raphson

Now we describe the second iterative step of our estimation procedure: assume an estimate  $\hat{\boldsymbol{\beta}}$  of  $\boldsymbol{\beta}$  is available and estimate the variance components  $\{\mathbf{G}, \sigma^2\}$ . Setting  $\boldsymbol{\beta} = \hat{\boldsymbol{\beta}}$ , the log-likelihood becomes

$$l(\hat{\boldsymbol{\beta}}, \mathbf{V}) = -\frac{1}{2} \log(|\mathbf{V}|) - \frac{1}{2} (\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}})' \mathbf{V}^{-1} (\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}). \quad (9)$$

The maximum likelihood estimates for  $\{\mathbf{G}, \sigma^2\}$  can be obtained by maximizing (9). A major criticism of these maximum likelihood estimates is their bias caused by the loss in the degrees of freedom from the estimation of  $\boldsymbol{\beta}$ . To correct this, Harville (1974) introduced a restricted maximum likelihood approach in which a correction term  $-\frac{1}{2} \log |\mathbf{X}' \mathbf{V}^{-1} \mathbf{X}|$  is added to (9). However, empirical results seem to suggest that there is no major difference between the two

likelihood approaches. In our work, we use maximum likelihood and the maximization of (9) is achieved via the Newton-Raphson (NR) algorithm (Lindstrom and Bates, 1988).

To further simplify our maximization problem, we replace  $\sigma^2$  in (9) with its maximum likelihood estimate:

$$\hat{\sigma}^2 = \frac{1}{n} \mathbf{r}' \mathbf{V}^{*-1} \mathbf{r} \quad \text{with} \quad \mathbf{V}^* = \mathbf{Z} \mathbf{G} \mathbf{Z}' + \mathbf{I} \quad \text{and} \quad \mathbf{r} = \mathbf{Y} - \mathbf{X} \hat{\boldsymbol{\beta}}.$$

Then our estimate for  $\mathbf{G}$  is given as the minimizer of the following negative profile likelihood function

$$p(\mathbf{G}|\hat{\boldsymbol{\beta}}) = \log(|\mathbf{V}^*|) + n \log(\mathbf{r}' \mathbf{V}^{*-1} \mathbf{r}). \quad (10)$$

Let  $\mathbf{A} = \mathbf{Z}' \mathbf{V}^{*-1} \mathbf{r} (\mathbf{Z}' \mathbf{V}^{*-1} \mathbf{r})'$ ,  $\mathbf{B} = \mathbf{Z}' \mathbf{V}^{*-1} \mathbf{Z}$ . Then, as shown in Appendix A, the gradient and the Hessian matrix for the NR algorithm w.r.t.  $\text{vec}(\mathbf{G})$  are, respectively,

$$\frac{\partial p(\mathbf{G}|\hat{\boldsymbol{\beta}})}{\partial \text{vec}(\mathbf{G})} = \text{vec} \mathbf{B} - \frac{1}{\hat{\sigma}^2} \text{vec} \mathbf{A}$$

and

$$\frac{\partial^2 p(\mathbf{G}|\hat{\boldsymbol{\beta}})}{\partial \text{vec}(\mathbf{G}) \partial \text{vec}(\mathbf{G})'} = \frac{\partial \text{vec}(\mathbf{G}')}{\partial \text{vec}(\mathbf{G})} \frac{\partial^2 p(\mathbf{G}|\hat{\boldsymbol{\beta}})}{\partial \text{vec}(\mathbf{G}') \partial \text{vec}(\mathbf{G})'},$$

where

$$\frac{\partial^2 p(\mathbf{G}|\hat{\boldsymbol{\beta}})}{\partial \text{vec}(\mathbf{G}') \partial \text{vec}(\mathbf{G})'} = -\mathbf{B} \otimes \mathbf{B} - \frac{1}{n \hat{\sigma}^4} \text{vec}(\mathbf{A}) \text{vec}(\mathbf{A})' + \frac{1}{\hat{\sigma}^2} (\mathbf{A} \otimes \mathbf{B} + \mathbf{B} \otimes \mathbf{A}).$$

In particular, denote the diagonal elements of  $\mathbf{G}$  as  $\boldsymbol{\theta}$ , we have

$$\frac{\partial p(\mathbf{G}|\hat{\boldsymbol{\beta}})}{\partial \boldsymbol{\theta}} = \frac{\partial \text{vec}(\mathbf{G})}{\partial \boldsymbol{\theta}} \frac{\partial p(\mathbf{G}|\hat{\boldsymbol{\beta}})}{\partial \text{vec}(\mathbf{G})} = \text{diag}(\mathbf{B}) - \frac{1}{\hat{\sigma}^2} \text{diag}(\mathbf{A}) \quad (11)$$

and

$$\begin{aligned} \frac{\partial^2 p(\mathbf{G}|\hat{\boldsymbol{\beta}})}{\partial \boldsymbol{\theta}' \partial \boldsymbol{\theta}} &= \frac{\partial \text{vec}(\mathbf{G})}{\partial \boldsymbol{\theta}} \frac{\partial^2 p(\mathbf{G}|\hat{\boldsymbol{\beta}})}{\partial \text{vec}(\mathbf{G}) \partial \text{vec}(\mathbf{G})'} \left( \frac{\partial \text{vec}(\mathbf{G})}{\partial \boldsymbol{\theta}} \right)' \\ &= -\mathbf{B} \odot \mathbf{B} - \frac{1}{n \hat{\sigma}^4} (\mathbf{A} \odot \mathbf{A}) + \frac{2}{\hat{\sigma}^2} (\mathbf{A} \odot \mathbf{B}), \end{aligned} \quad (12)$$

where  $\odot$  represents the element-wise multiplication. With these expressions, the standard NR algorithm can be applied to maximize (10) and obtain estimates for  $\mathbf{V}$  and  $\mathbf{G}$ . We note Woodbury's identity  $\mathbf{V}^{*-1} = \mathbf{I} - \mathbf{Z}(\mathbf{Z}'\mathbf{Z} + \mathbf{G}^{-1})^{-1}\mathbf{Z}'$  can be applied to avoid heavy computations.

### 3.3. Combining the two algorithms

The above BCGD and NR algorithms can be applied iteratively to obtain estimates for both the fixed and random components parameters  $\beta$ ,  $\mathbf{V}$  and  $\mathbf{G}$ . The combined algorithm begins with initial estimates  $\hat{\beta}^{(0)}$  and  $\hat{\mathbf{V}}^{(0)}$  for  $\beta$  and  $\mathbf{V}$ , respectively, and iterate, for  $t = 0, 1, \dots$ , the following two steps until convergence:

1. *Fixed Effects Estimation:* Obtain the next iterative estimate  $\hat{\beta}^{(t+1)}$  with the BCGD Algorithm described in Section 3.1.
2. *Random Effects Estimation:* Obtain  $\hat{\mathbf{V}}^{(t+1)}$  and  $\hat{\mathbf{G}}^{(t+1)}$  by minimizing (10) with the NR algorithm described in Section 3.2.

Convergence of this combined algorithm can be determined by monitoring the successive changes of the iterative fitted value of  $\mathbf{Y}$ . In our implementation we set the initial estimates as  $\hat{\beta}^{(0)} = 0$  and  $\hat{\mathbf{V}}^{(0)} = \mathbf{I}$ .

Denote the final estimates as  $\hat{\beta}$ ,  $\hat{\mathbf{V}}$  and  $\hat{\mathbf{G}}$ . When the above iteration finishes, the estimated conditional expectation of  $\mathbf{u}$  given  $\beta = \hat{\beta}$  can be calculated as

$$\hat{\mathbf{u}} = (\mathbf{Z}'\mathbf{Z} + \hat{\mathbf{G}}^{-1})^{-1}\mathbf{Z}'(\mathbf{Y} - \mathbf{X}\hat{\beta}), \quad (13)$$

while the estimate for  $E[\mathbf{Y}|\mathbf{X}, \mathbf{u}]$  is given by  $\hat{\mathbf{Y}} = \mathbf{X}\hat{\beta} + \mathbf{Z}\hat{\mathbf{u}}$ . Also,  $\hat{\mathbf{Y}}$  can be used as the prediction of  $\mathbf{Y}$ .

## 4. Selection of model complexity

This section develops a method for selecting the model complexity of (1). That is, to select the adaptive group lasso parameters  $(\lambda, \nu)$  and to determine which random effects  $u_q$ 's should enter the final model. Note that the selection of the fixed effects  $f_p$ 's (which is equivalent to  $\beta_p$ 's) is achieved by varying the values of  $(\lambda, \nu)$ .

### 4.1. Model selection criterion

We adopt the Bayesian information criterion (BIC) to select the model complexity. To be specific, we define the best fitting model as the one that minimizes the following objection function:

$$n \log(\hat{\sigma}^2) + \frac{1}{\hat{\sigma}^2} \|\mathbf{Y} - \mathbf{X}\hat{\beta} - \mathbf{Z}\hat{\mathbf{u}}\|^2 + \text{df}_{\hat{\mathbf{Y}}} \cdot \log(n), \quad (14)$$

where  $\text{df}_{\hat{\mathbf{Y}}}$  is the effective degrees of freedom of the fitted model. Computable expression for  $\text{df}_{\hat{\mathbf{Y}}}$  is given in the next subsection. Notice that the first two terms in (14) together form the conditional likelihood of  $\mathbf{Y}$  given  $\hat{\beta}$  and  $\hat{\mathbf{u}}$  (e.g., Ruppert *et al.*, 2003).



#### 4.2. Calculating the degrees of freedom

This subsection derives a computable expression for  $df_{\hat{\mathbf{Y}}}$  so that (14) can be used in practice.

For simplicity, suppose for the moment that  $\mathbf{V}$  is known or can be independently estimated. The degrees of freedom of any fitted model  $df_{\hat{\mathbf{Y}}}$  is calculated as the sum of the degrees of freedom contributed by the fixed and random effects parameters. If we denote the degrees of freedom of the fixed and random effects as, respectively,  $df_{\hat{\boldsymbol{\beta}}}$  and  $df_{\hat{\mathbf{u}}}$ , we have  $df_{\hat{\mathbf{Y}}} = df_{\hat{\boldsymbol{\beta}}} + df_{\hat{\mathbf{u}}}$ .

Following Shen and Ye (2002) and Ruppert *et al.* (2003), we define  $df_{\hat{\boldsymbol{\beta}}}$  and  $df_{\hat{\mathbf{u}}}$  as, respectively,

$$df_{\hat{\boldsymbol{\beta}}} = E\left\{\text{tr}\left(\frac{\partial \mathbf{X}\hat{\boldsymbol{\beta}}}{\partial \mathbf{Y}}\right)\right\} = 1 + E\left\{\text{tr}\left(\sum_{p=1}^P \frac{\partial \hat{\boldsymbol{\beta}}_p}{\partial \mathbf{Y}} \mathbf{X}'_p\right)\right\}, \quad (15)$$

and

$$df_{\hat{\mathbf{u}}} = E\left\{\text{tr}\left(\frac{\partial \mathbf{Z}\hat{\mathbf{u}}}{\partial \mathbf{Y}}\right)\right\} = E\left[\text{tr}\left\{\mathbf{Z}(\mathbf{Z}'\mathbf{Z} + \mathbf{G}^{-1})^{-1}\mathbf{Z}'\left(\mathbf{I} - \frac{1}{n}\mathbf{J} - \sum_{p=1}^P \frac{\partial \hat{\boldsymbol{\beta}}_p}{\partial \mathbf{Y}} \mathbf{X}'_p\right)\right\}\right], \quad (16)$$

where  $\mathbf{J}$  is a matrix of ones. In these expressions the only unknown is  $\frac{\partial \hat{\boldsymbol{\beta}}_p}{\partial \mathbf{Y}}$ , and it can be calculated as follows.

By the KKT conditions of (5), we have, for  $p = 1, \dots, P$ ,

$$-\mathbf{X}'_p \mathbf{V}^{-1}(\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}) + \lambda w_p \frac{\hat{\boldsymbol{\beta}}_p}{\|\hat{\boldsymbol{\beta}}_p\|} = \mathbf{0} \quad \text{if } \hat{\boldsymbol{\beta}}_p \neq \mathbf{0},$$

which gives

$$\frac{\partial \hat{\boldsymbol{\beta}}_p}{\partial \mathbf{Y}} \left\{ \mathbf{X}'_p \mathbf{V}^{-1} \mathbf{X}_p + \frac{\lambda w_p}{\|\hat{\boldsymbol{\beta}}_p\|} \left( \mathbf{I} - \frac{\hat{\boldsymbol{\beta}}_p \hat{\boldsymbol{\beta}}_p'}{\|\hat{\boldsymbol{\beta}}_p\|^2} \right) \right\} + \sum_{j \neq p} \frac{\partial \hat{\boldsymbol{\beta}}_j}{\partial \mathbf{Y}} \mathbf{X}'_j \mathbf{V}^{-1} \mathbf{X}_p = \left( \mathbf{I} - \frac{1}{n} \mathbf{J} \right) \mathbf{V}^{-1} \mathbf{X}_p.$$

The above can be expressed in a more compact manner as follows:

$$\left( \frac{\partial \hat{\boldsymbol{\beta}}_1}{\partial \mathbf{Y}}, \dots, \frac{\partial \hat{\boldsymbol{\beta}}_P}{\partial \mathbf{Y}} \right) \mathbf{C} = \left( \mathbf{D}_1, \dots, \mathbf{D}_P \right) = \mathbf{D},$$

where  $\mathbf{C} = \mathbf{X}' \mathbf{V}^{-1} \mathbf{X} + \mathbf{E}$ ,  $\mathbf{E} = \text{diag}(\mathbf{E}_1, \dots, \mathbf{E}_P)$ ,

$$\mathbf{D}_p = \begin{cases} \left( \mathbf{I} - \frac{1}{n} \mathbf{J} \right) \mathbf{V}^{-1} \mathbf{X}_p, & \text{if } \hat{\boldsymbol{\beta}}_p \neq \mathbf{0} \\ \mathbf{0}, & \text{if } \hat{\boldsymbol{\beta}}_p = \mathbf{0} \end{cases}$$

and

$$\mathbf{E}_p = \begin{cases} \frac{\lambda w_p}{\|\hat{\boldsymbol{\beta}}_p\|} \left( \mathbf{I} - \frac{\hat{\boldsymbol{\beta}}_p \hat{\boldsymbol{\beta}}_p'}{\|\hat{\boldsymbol{\beta}}_p\|^2} \right), & \text{if } \hat{\boldsymbol{\beta}}_p \neq \mathbf{0} \\ \mathbf{0}, & \text{if } \hat{\boldsymbol{\beta}}_p = \mathbf{0} \end{cases}.$$

Therefore we have

$$\left[ \frac{\partial \hat{\beta}_1}{\partial \mathbf{Y}}, \dots, \frac{\partial \hat{\beta}_P}{\partial \mathbf{Y}} \right] = \mathbf{D}\mathbf{C}^{-1}.$$

We note that  $\mathbf{X}'\mathbf{V}^{-1}\mathbf{X}$  is typically not of full rank, hence  $\mathbf{C}$  may not be invertible in some cases (e.g., when  $\lambda = 0$ ). However,  $\mathbf{X}\mathbf{C}^{-1}\mathbf{D}'$  is invariant w.r.t. the choice of the pseudo inverse of  $\mathbf{C}$ . Hence, we can estimate (15) and (16) by

$$\hat{\text{df}}_{\beta} = 1 + \text{tr} \left( \sum_{p=1}^P \frac{\partial \hat{\beta}_p}{\partial \mathbf{Y}} \mathbf{X}'_p \right), \quad (17)$$

and

$$\hat{\text{df}}_{\hat{\mathbf{u}}} = \left\{ \mathbf{Z}(\mathbf{Z}'\mathbf{Z} + \hat{\mathbf{G}}^{-1})^{-1} \mathbf{Z}' \left( \mathbf{I} - \frac{1}{n} \mathbf{J} - \sum_{p=1}^P \frac{\partial \hat{\beta}_p}{\partial \mathbf{Y}} \mathbf{X}'_p \right) \right\}, \quad (18)$$

respectively. And the estimated degrees of freedom for the model equals  $\hat{\text{df}}_{\hat{\mathbf{Y}}} = \hat{\text{df}}_{\beta} + \hat{\text{df}}_{\hat{\mathbf{u}}}$ . In practice we replace  $\mathbf{V}$  by the  $\hat{\mathbf{V}}$  obtained using the algorithm summarized in Section 3.3.

#### 4.3. Practical minimization of (14)

In practice a global minimization of the selection criterion (14) is difficult. Here we suggest using the following procedure to approximate the minimizer of (14). This procedure is relatively fast, and produces excellent numerical results in all our experimental work. The steps are:

1. *Initialization:* With the full model (i.e., with all fixed and random effects), apply the ordinary group lasso to obtain initial parameter estimates. That is, the initial estimates are chosen as the minimizers of (5) with  $w_p = 1$  for all  $p$ . The tuning parameter of this ordinary group lasso can be chosen for example by BIC.
2. *Fixed Effects Selection:* Using the estimates obtained in Step 1 as initial estimates and weights (6), apply the proposed parameter estimation method summarized in Section 3.3 to the full model. The adaptive group lasso parameters  $\lambda$  and  $\nu$  are chosen by the BIC criterion (14) through a two-dimensional grid search.
3. *Random Effects Selection:* Perform the following backward selection strategy to the fitted model obtained in Step 2. First from this model remove the random effect that has the smallest coefficient magnitude and re-calculate its BIC (14) value, denoted as  $\text{BIC}_1$ . If  $\text{BIC}_1$  is larger than the BIC value of the model from Step 2, then the model from Step 2 is taken as the final model. Otherwise, further remove the random effect that has the second smallest magnitude and re-calculate its corresponding BIC value; denoted it as  $\text{BIC}_2$ . If  $\text{BIC}_2 > \text{BIC}_1$ , then the earlier model with one random effect removed is taken as the final model. Otherwise, repeat

this process by removing more random effects until there is no decrease in the BIC values.

One could iterate Steps 2 and 3 above, but our experience suggests that the improvement, if any, is minimal.

## 5. Theoretical results

This section establishes the consistency properties of the proposed adaptive group lasso estimator, as well as the BIC criterion (14).

### 5.1. Consistency property of adaptive group lasso

Denote  $A_0 \equiv \{p : \|f_p\| = 0, p = 1, \dots, P\}$  and  $A_T \equiv \{p : \|f_p\| > 0, p = 1, \dots, P\}$ . That is,  $A_0$  is the set of all non-significant  $f_p$ 's, while  $A_T$  is the set of all significant  $f_p$ 's. Let  $\mathbf{M}$  be any matrix in  $\mathbb{R}^{m \times n}$ . Define the matrix norm of  $\mathbf{M}$  as  $\|\mathbf{M}\| = (\max_{\mathbf{x} \in \mathbb{R}^n, \|\mathbf{x}\|=1} \mathbf{x}'\mathbf{M}'\mathbf{M}\mathbf{x})^{1/2}$ . An important property of this matrix norm is that, for any  $\mathbf{x} \in \mathbb{R}^n$ ,  $\|\mathbf{M}\mathbf{x}\| \leq \|\mathbf{M}\|\|\mathbf{x}\|$ . We consider the following regularity conditions (A.1-7):

- (A.1) The number of elements in  $A_T$ , denoted as  $|A_T|$ , does not depend on  $n$  and  $P$ .
- (A.2) For the knots  $(t_i, i = 0, \dots, \tau + 1)$ , let  $h_i = t_{i+1} - t_i$  and  $h = \max_i h_i$ . We assume that  $h/\min_i h_i = O(1)$  and  $\max_i |h_{i+1} - h_i| = o(m^{-1})$ .
- (A.3) Without loss of generality, we assume that the design points  $X_{ip} \in [0, 1]$ . Moreover, we assume that  $X_{ip}$  has a nonzero continuous density function.
- (A.4) Let  $\mathcal{F}$  be the collection of functions satisfying Lipschitz condition:

$$|f^{(k)}(s) - f^{(k)}(t)| \leq C|s - t|^\alpha, \quad s, t \in [0, 1], \text{ for some } C > 0$$

where  $k$  is non-negative integer,  $0 < \alpha \leq 1$  and  $0.5 < k + \alpha = d$  is called the order of smoothness. For  $p = 1, \dots, P$ ,  $f_p \in \mathcal{F}$  and  $\text{E}f_p(X_{ip}) = 0$ .

- (A.5) The initial estimate  $\tilde{\boldsymbol{\beta}}$  of  $\boldsymbol{\beta}$  satisfies the following conditions as  $n \rightarrow \infty$ :

$$(a) \quad r \max_{p \in A_0} \|\tilde{\boldsymbol{\beta}}_p\| = O_p(1) \text{ with } r \rightarrow \infty.$$

$$(b) \quad P(\min_{p \in A_T} \|\tilde{\boldsymbol{\beta}}_p\| > \nu^{-1}) \rightarrow 1 \text{ for any } \nu > 0.$$

- (A.6) As  $n \rightarrow \infty$ :

$$(a) \quad m \rightarrow \infty, \frac{m}{n^{1/2}} \rightarrow 0, \frac{m \log(2mP)}{n} \rightarrow 0.$$

$$(b) \quad \frac{n}{\lambda r m^{d+1/2}} \rightarrow 0, \frac{\sqrt{n \log(2mP)}}{\lambda r} \rightarrow 0.$$

- (A.7) The norm of the covariance matrix  $\|\mathbf{V}\|$  is finite.

We need more notation to proceed. For any  $A \subset \{1, \dots, P\}$ , denote  $\mathbf{X}_A$  and  $\boldsymbol{\beta}_A$ , respectively, as the design matrix and the subset of  $\boldsymbol{\beta}$  corresponding to the

model with fixed effects  $f_p$ 's indexed by  $A$ . Then the generalized least squares estimate of  $\beta_A$  based on  $\mathbf{X}_A$  is  $\hat{\mathbf{b}}_A = (\hat{\mathbf{b}}'_{A,p}, p = 1, \dots, P)'$ , where  $(\hat{\mathbf{b}}'_{A,p}, p \in A)' = (\mathbf{X}'_A \mathbf{V}^{-1} \mathbf{X}_A)^{-1} \mathbf{X}'_A \mathbf{V}^{-1} \mathbf{Y}$  and  $(\hat{\mathbf{b}}'_{A,p}, p \notin A)' = \mathbf{0}$ . Without loss of generality, we assume  $\mu = 0$ .

The following theorem shows that our adaptive group lasso estimate of  $\beta$  is not only capable of identifying the correct variables, but is also capable of recovering the predictive performance of the generalized least squares estimate based on the correct variables. The proof can be found in Appendix C.

**Theorem 1** ( $\mathbf{G}$  known). *Under (A.1)-(A.7), as  $n \rightarrow \infty$ ,*

1.  $P(\text{sign}(\|\hat{\beta}_p\|) = \text{sign}(\|\beta_p\|), p = 1, \dots, P) \rightarrow 1$ .
2.  $P(\hat{\beta} = \hat{\mathbf{b}}_{A_T}) \rightarrow 1$ .
3.  $\sum_{p=1}^P \|\hat{\beta}_p - \beta_p\|^2 = O_p\left(\frac{m^2}{n}\right) + O_p\left(\frac{1}{m^{2d-1}}\right) + O_p\left(\frac{m^2 \lambda^2}{n^2}\right) o_p(1)$ .
4.  $\sum_{p=1}^P \|\hat{f}_p - f_p\|^2 = O_p\left(\frac{m}{n}\right) + O_p\left(\frac{1}{m^{2d}}\right) + O_p\left(\frac{m \lambda^2}{n^2}\right) o_p(1)$ .

**Corollary 1** ( $\mathbf{G}$  unknown). *Under the conditions of Theorem 1, except that  $\mathbf{G}$  is replaced by a consistent estimate  $\hat{\mathbf{G}}$ , Theorem 1 continue to hold.*

*Proof.* The proof is essentially the same as that for Theorem 1 after applying Slutsky's theorem:  $h(\hat{\mathbf{G}}) \rightarrow h(\mathbf{G})$  in probability for any real-value continuous function  $h$ .  $\square$

## 5.2. Consistency property of Bayesian information criterion

First assume  $\mathbf{G}$  is known. Then the maximum likelihood estimate of  $\sigma^2$  based on model  $A$  is  $\hat{\sigma}_A^2 \equiv \frac{1}{n}(\mathbf{Y} - \mathbf{X}\hat{\mathbf{b}}_A)' \mathbf{V}^{*-1}(\mathbf{Y} - \mathbf{X}\hat{\mathbf{b}}_A)$ . To stress the dependence of the parameter estimates on  $(\lambda, \nu)$ , denote  $\hat{\beta}_{\lambda, \nu}$  and  $\hat{\mathbf{u}}_{\lambda, \nu}$  as the adaptive grouped lasso estimates with weights (6) and tuning parameters  $(\lambda, \nu)$ . Also let  $\hat{\sigma}_{\lambda, \nu}^2 \equiv \frac{1}{n}(\mathbf{Y} - \mathbf{X}\hat{\beta}_{\lambda, \nu})' \mathbf{V}^{*-1}(\mathbf{Y} - \mathbf{X}\hat{\beta}_{\lambda, \nu})$  be the corresponding estimate of  $\sigma^2$ . With this, the Bayesian information criterion (14) can be re-written as:

$$\text{BIC}(\lambda, \nu) = \log(\hat{\sigma}_{\lambda, \nu}^2) + \frac{1}{n\hat{\sigma}_{\lambda, \nu}^2} \|\mathbf{Y} - \mathbf{X}\hat{\beta}_{\lambda, \nu} - \mathbf{Z}\hat{\mathbf{u}}_{\lambda, \nu}\|^2 + \hat{\text{df}}_{\hat{\mathbf{Y}}}(\lambda, \nu) \frac{\log(n)}{n}, \quad (19)$$

where  $\hat{A}_{\lambda, \nu} \equiv \{p : \|\hat{\beta}_{\lambda, \nu, p}\| > 0\}$ . We consider the following technical conditions:

- (B.1) For any  $A \in \mathcal{A}$ , there exists  $\sigma_A^2 > 0$  such that  $\hat{\sigma}_A^2 \rightarrow \sigma_A^2$  in probability.
- (B.2) For any  $A$  with  $A \cap A_T \neq A_T$ ,  $\sigma_A^2 > \sigma_{A_T}^2$ .
- (B.3) As  $n \rightarrow \infty$ ,  $\lambda/n \rightarrow 0$ .

Appendix C establishes the following theorem.

**Theorem 2** ( $G$  known). Under (A.1)-(A.7) and (B.1)-(B.3). Assume that (A.5) and (A.6) hold for some sequences  $\{\lambda_n^*\}$  and  $\{\nu_n^*\}$ . Let

$$(\hat{\lambda}, \hat{\nu}) = \underset{(\lambda, \nu)}{\operatorname{arg\,min}} \operatorname{BIC}(\lambda, \nu).$$

Then

$$P(\hat{\beta}_{\hat{\lambda}, \hat{\nu}} = \hat{b}_{A_T}) \rightarrow 1 \quad \text{as } n \rightarrow \infty.$$

As similar to Corollary 1, applying Slutsky’s theorem leads to the following Corollary.

**Corollary 2** ( $G$  unknown). Under the conditions of Theorem 2, except that  $G$  is replaced by a consistent estimate  $\hat{G}$ , Theorem 2 continues to hold.

**6. Simulations**

In this section, we investigate the empirical performance of our proposed methodology via numerical experiments. The covariates  $X_{ip}$ ’s were generated from the uniform(0, 1) distribution. The test functions (i.e., fixed effects) used in the experiments are listed in Table 1; these functions have been used by Meier *et al.* (2009). Note that the functions listed in Table 1 are not centered, while in all our numerical work, they were, however, all centered.

We tested our method with three different types of models:

**Model 1 (categorical variables as mixed-effects):**

$$\begin{aligned} Y_i = & f_1(X_{i1}) + f_2(X_{i2}) + f_3(X_{i3}) + f_4(X_{i4}) + \sum_{p=5}^{20} f_p(X_{ip}) \\ & + \sum_{k=1}^5 u_{1k} I(Z_{i1} = k) + \sum_{k=1}^5 u_{2k} I(Z_{i2} = k) \\ & + \sum_{k=1}^3 u_{3k} I(Z_{i3} = k) + \sum_{k=1}^3 u_{4k} I(Z_{i4} = k) + \epsilon_i, \end{aligned}$$

where  $u_{qk}$ ’s are i.i.d.  $N(0, \theta_q^2 \sigma^2)$  with  $\theta_1 = 3$ ,  $\theta_2 = 4$ , and  $\theta_3 = \theta_4 = 0$ , and  $\epsilon_i$ ’s are i.i.d.  $N(0, \sigma^2)$ . The  $Z_{iq}$ ’s are generated from the discrete uniform[1, . . . , 5] distribution for  $q = 1, 2$  and discrete uniform[1, 2, 3] for  $q = 3, 4$ .

TABLE 1  
Test functions used in the simulation, all were centered in the simulation

test function	
$f_1$	$6x$
$f_2$	$5(2x - 1)^2$
$f_3$	$4 \sin(2\pi x) / \{2 - \sin(2\pi x)\}$
$f_4$	$3\{0.1 \sin(2\pi x) + 0.2 \cos(2\pi x) + 0.3 \sin^2(2\pi x) + 0.4 \cos^2(2\pi x) + 0.5 \sin^3(2\pi x)\}$
$f_p, p \geq 5$	0

**Model 2 (random intercept):**

$$Y_{ik} = u_i + f_1(X_{ik1}) + f_2(X_{ik2}) + f_3(X_{ik3}) + f_4(X_{ik4}) \\ + \sum_{p=5}^{20} f_p(X_{ikp}) + \epsilon_{ik},$$

where  $u_i$ 's are i.i.d.  $N(0, \theta^2 \sigma^2)$  with  $\theta = 3$  and  $\epsilon_{ik}$ 's are i.i.d.  $N(0, \sigma^2)$ . The response  $Y_{ik}$  represents the  $k$ th measurement of the  $i$ th subject, where  $1 \leq i \leq N$ ,  $1 \leq k \leq n_i$  and  $\sum_{i=1}^N n_i = n$ . Here  $N$  is the total number of subjects and  $n_i$  is the number of measurements taken from the  $i$ th subject.

**Model 3 (random intercepts and trends):**

$$Y_{ik} = a_i + b_{i1}X_{ik1} + f_3(X_{ik1}) + b_{i2}X_{ik2} + f_4(X_{ik2}) \\ + \sum_{p=3}^6 \{b_{ip}X_{ikp} + f_{p+2}(X_{ikp})\} + \epsilon_{ik},$$

where  $a_i \sim N(0, \theta_1^2 \sigma^2)$ ,  $b_{i1} \sim N(0, \theta_2^2 \sigma^2)$ ,  $b_{i2} \sim N(0, \theta_3^2 \sigma^2)$  with  $\theta_i = i + 1$  and  $b_{ip} = 0$  for  $p \geq 3$ , and  $\epsilon_{ik}$ 's are i.i.d.  $N(0, \sigma^2)$ . As above,  $Y_{ik}$  represents the  $k$ th measurement of the  $i$ th subject, where  $1 \leq i \leq N$ ,  $1 \leq k \leq n_i$  and  $\sum_{i=1}^N n_i = n$ .

We tested the proposed method with different combinations of  $n$ ,  $N$  and signal-to-noise ratio (SNR, defined as  $\text{SNR} = \text{Var}\{E(Y|\mathbf{X})\}^{\frac{1}{2}}/\sigma$ ). The number of Monte-Carlo runs tested for each experimental configuration was 200. We use the following mean squared error (MSE) to measure the estimation accuracy for the additive components:  $\text{MSE} = \sum_{p=1}^P \|f_p - \hat{f}_p\|^2$ .

For each simulated data set, various estimates were obtained by three different estimation methods:

- Method I: applying the proposed method to the full data set.
- Method II: applying the ordinary group lasso method to the full data set (i.e., with both fixed and random effects covariates). No random effects selection is performed as all the random effects are always included in the final model. The tuning parameter  $\lambda$  for the group lasso selection is chosen by BIC. This method is implemented in a similar fashion as the approach described in Section 3.3, except no random effects selection is done.
- Method III: this method is used here as an ‘‘oracle method’’ for benchmark comparison, as only the significant variables are considered (i.e., assuming the true model is known). It iterates between the following two steps: use generalized least squares to estimate the fixed effects parameters and use the NR algorithm to estimate the random effects parameters.

For each estimated model, we counted the number of non-zero  $\hat{f}_p$ , and the number of non-zero random effects, intercepts and/or trends. We also calculated the above-defined MSE, and recorded the parameter estimates for the random components. Lastly, we noted if the estimated model is exactly the same as the true model, or if it is a superset of the true model. These results are summarized in Tables 2 to 7.

TABLE 2  
Simulation results for Model 1. The true numbers of non-zero  $f_p$ 's and random effects are 4 and 2 respectively

SNR = 3		number of non-zero $\hat{f}_p$				MSE		number of random effects			number of times true model is	
configuration	method	4	5	6	7+	mean	S.E.	2	3	4	included	selected
$n = 128$ $m = 7$	I	193	5	2	0	0.3748	0.0067	161	36	3	200	155
	II	39	48	40	73	0.6884	0.0136	—	—	—	—	—
	III	200	0	0	0	0.3753	0.0069	—	—	—	—	—
$n = 256$ $m = 8$	I	200	0	0	0	0.2467	0.0037	167	30	3	200	167
	II	118	52	18	12	0.4284	0.0071	—	—	—	—	—
	III	200	0	0	0	0.2082	0.0033	—	—	—	—	—
$n = 512$ $m = 9$	I	200	0	0	0	0.1351	0.0019	186	12	2	200	186
	II	156	41	3	0	0.2152	0.0035	—	—	—	—	—
	III	200	0	0	0	0.1105	0.0014	—	—	—	—	—

SNR = 4		number of non-zero $\hat{f}_p$				MSE		number of random effects			number of times true model is	
configuration	method	4	5	6	7+	mean	S.E.	2	3	4	included	selected
$n = 128$ $m = 7$	I	200	0	0	0	0.2605	0.0038	161	34	5	200	161
	II	58	49	36	57	0.5021	0.009	—	—	—	—	—
	III	200	0	0	0	0.2461	0.004	—	—	—	—	—
$n = 256$ $m = 8$	I	200	0	0	0	0.1835	0.0025	170	29	1	200	170
	II	123	53	20	4	0.3073	0.0045	—	—	—	—	—
	III	200	0	0	0	0.1501	0.0019	—	—	—	—	—
$n = 512$ $m = 9$	I	200	0	0	0	0.0926	0.0011	182	16	2	200	182
	II	154	40	4	2	0.1434	0.0018	—	—	—	—	—
	III	200	0	0	0	0.0747	9e-04	—	—	—	—	—

TABLE 3  
Estimates of random components for Model 1

SNR = 3		$\hat{\theta}_1$ ( $\theta_1 = 3$ )			$\hat{\theta}_2$ ( $\theta_2 = 4$ )			$\hat{\sigma}$ ( $\sigma \approx 1$ )		
configuration	method	median	mean	s.e.	median	mean	s.e.	median	mean	s.e.
$n = 128$ $m = 7$	I	2.874	2.888	0.071	3.621	3.767	0.092	1.029	1.028	0.006
	II	2.579	2.709	0.089	3.271	3.556	0.121	1.131	1.132	0.009
	III	2.936	2.983	0.075	3.75	3.89	0.095	0.994	0.996	0.006
$n = 256$ $m = 8$	I	2.596	2.623	0.063	3.429	3.523	0.087	1.093	1.093	0.005
	II	2.393	2.45	0.06	3.193	3.299	0.084	1.176	1.173	0.005
	III	2.684	2.723	0.066	3.558	3.655	0.09	1.053	1.055	0.005
$n = 512$ $m = 9$	I	2.597	2.645	0.069	3.361	3.414	0.09	1.101	1.1	0.003
	II	2.499	2.55	0.066	3.252	3.295	0.088	1.139	1.138	0.003
	III	2.656	2.707	0.07	3.446	3.493	0.092	1.075	1.075	0.003

SNR = 4		$\hat{\theta}_1$ ( $\theta_1 = 3$ )			$\hat{\theta}_2$ ( $\theta_2 = 4$ )			$\hat{\sigma}$ ( $\sigma \approx 0.8$ )		
configuration	method	median	mean	s.e.	median	mean	s.e.	median	mean	s.e.
$n = 128$ $m = 7$	I	2.717	2.735	0.07	3.477	3.727	0.122	0.811	0.818	0.005
	II	2.423	2.505	0.08	3.101	3.551	0.19	0.911	0.917	0.007
	III	2.82	2.835	0.072	3.606	3.764	0.093	0.778	0.789	0.005
$n = 256$ $m = 8$	I	2.601	2.631	0.071	3.481	3.61	0.082	0.839	0.84	0.004
	II	2.432	2.437	0.066	3.233	3.343	0.075	0.91	0.907	0.004
	III	2.689	2.735	0.074	3.627	3.753	0.085	0.808	0.809	0.003
$n = 512$ $m = 9$	I	2.537	2.614	0.064	3.257	3.335	0.078	0.837	0.835	0.003
	II	2.477	2.53	0.062	3.153	3.225	0.075	0.868	0.863	0.003
	III	2.606	2.677	0.066	3.319	3.415	0.08	0.817	0.816	0.002

TABLE 4  
Simulation results for Model 2. The true number of  $f_p$ 's is 4, and  $N_I$  denotes "number of times the random intercept is selected"

SNR = 3		number of non-zero $\hat{f}_p$				MSE			number of times true model is	
configuration	method	4	5	6	7+	mean	S.E.	$N_I$	included	selected
$n = 128$ $m = 7$ $N = 16$	I	195	5	0	0	0.3944	0.0068	200	200	195
	II	34	65	46	55	0.7134	0.0138	—	—	—
	III	200	0	0	0	0.4202	0.0077	—	—	—
$n = 256$ $m = 8$ $N = 16$	I	200	0	0	0	0.2519	0.0041	200	200	200
	II	109	66	18	7	0.4376	0.0074	—	—	—
	III	200	0	0	0	0.2166	0.0033	—	—	—
$n = 512$ $m = 9$ $N = 16$	I	200	0	0	0	0.1344	0.0018	200	200	200
	II	168	32	0	0	0.2203	0.0031	—	—	—
	III	200	0	0	0	0.112	0.0016	—	—	—
$n = 256$ $m = 8$ $N = 32$	I	200	0	0	0	0.2659	0.0036	200	200	200
	II	103	65	26	6	0.4744	0.0074	—	—	—
	III	200	0	0	0	0.2341	0.0035	—	—	—
$n = 512$ $m = 9$ $N = 32$	I	200	0	0	0	0.14	0.002	200	200	200
	II	157	37	5	1	0.226	0.0036	—	—	—
	III	200	0	0	0	0.1148	0.0015	—	—	—

SNR = 4		number of non-zero $\hat{f}_p$				MSE			number of times true model is	
configuration	method	4	5	6	7+	mean	S.E.	$N_I$	included	selected
$n = 128$ $m = 7$ $N = 16$	I	197	3	0	0	0.2763	0.0041	200	200	197
	II	40	51	42	67	0.5337	0.009	—	—	—
	III	200	0	0	0	0.265	0.0046	—	—	—
$n = 256$ $m = 8$ $N = 16$	I	200	0	0	0	0.1783	0.0023	200	200	200
	II	131	50	17	2	0.3039	0.0046	—	—	—
	III	200	0	0	0	0.1503	0.0018	—	—	—
$n = 512$ $m = 9$ $N = 16$	I	200	0	0	0	0.0959	0.0012	200	200	200
	II	160	36	3	1	0.1484	0.002	—	—	—
	III	200	0	0	0	0.0765	0.001	—	—	—
$n = 256$ $m = 8$ $N = 32$	I	200	0	0	0	0.1859	0.0026	200	200	200
	II	115	52	27	6	0.3191	0.0049	—	—	—
	III	200	0	0	0	0.1578	0.0021	—	—	—
$n = 512$ $m = 9$ $N = 32$	I	200	0	0	0	0.0945	0.0012	200	200	200
	II	171	27	2	0	0.1505	0.0021	—	—	—
	III	200	0	0	0	0.0764	0.001	—	—	—

The simulation results seem to suggest that our proposed method can often correctly select the correct significant fixed effects components  $f_p$ 's. It is also capable of selecting the significant random effects. The empirical performances of the proposed method improves as the sample size and/or SNR increase. The simulation results also show that our method outperforms the ordinary group lasso in terms of fixed effects selection.

When estimating the variance components (e.g.,  $\theta_1$  and  $\theta_2$  in Model 1), a bias pattern seems to exist: the bias increases as  $n$  increases. This is not too surprising, as it is known that the maximum likelihood estimates of the variance components are biased, as they do not account for the loss in degrees of freedom



TABLE 5  
Estimates of random components for Model 2

SNR = 3		$\hat{\theta} (\theta = 3)$			$\hat{\sigma} (\sigma \approx 1)$		
configuration	method	median	mean	s.e.	median	mean	s.e.
$n = 128$	I	3.085	3.171	0.046	1.044	1.036	0.007
$m = 7$	II	2.763	2.84	0.044	1.163	1.16	0.008
$N = 16$	III	3.169	3.267	0.048	1.007	1.004	0.006
$n = 256$	I	2.843	2.828	0.042	1.093	1.099	0.005
$m = 8$	II	2.662	2.631	0.039	1.174	1.182	0.005
$N = 16$	III	2.943	2.934	0.044	1.054	1.059	0.005
$n = 512$	I	2.884	2.891	0.037	1.098	1.095	0.003
$m = 9$	II	2.786	2.782	0.036	1.14	1.138	0.003
$N = 16$	III	2.955	2.957	0.038	1.073	1.071	0.003
$n = 256$	I	2.948	2.98	0.031	1.102	1.098	0.005
$m = 8$	II	2.719	2.748	0.028	1.197	1.19	0.005
$N = 32$	III	3.076	3.091	0.032	1.063	1.059	0.005
$n = 512$	I	2.954	2.961	0.026	1.097	1.097	0.003
$m = 9$	II	2.846	2.855	0.026	1.138	1.139	0.004
$N = 32$	III	3.026	3.028	0.027	1.071	1.073	0.003

SNR = 4		$\hat{\theta} (\theta = 3)$			$\hat{\sigma} (\sigma \approx 0.8)$		
configuration	method	median	mean	s.e.	median	mean	s.e.
$n = 128$	I	2.948	3.478	0.179	0.805	0.802	0.005
$m = 7$	II	2.603	4.357	0.369	0.918	0.912	0.007
$N = 16$	III	3.035	3.074	0.044	0.778	0.777	0.005
$n = 256$	I	2.819	2.811	0.039	0.844	0.836	0.004
$m = 8$	II	2.576	2.585	0.036	0.913	0.91	0.004
$N = 16$	III	2.929	2.926	0.04	0.81	0.804	0.004
$n = 512$	I	2.797	2.797	0.039	0.833	0.831	0.003
$m = 9$	II	2.666	2.693	0.038	0.864	0.863	0.003
$N = 16$	III	2.861	2.864	0.04	0.813	0.811	0.003
$n = 256$	I	2.917	2.912	0.029	0.834	0.835	0.004
$m = 8$	II	2.701	2.675	0.027	0.909	0.91	0.004
$N = 32$	III	3.037	3.025	0.03	0.805	0.803	0.004
$n = 512$	I	2.86	2.911	0.027	0.828	0.83	0.003
$m = 9$	II	2.726	2.794	0.026	0.864	0.865	0.003
$N = 32$	III	2.926	2.98	0.027	0.81	0.811	0.003

when estimating the fixed effects (Harville, 1974). However, for Model 1, we believe that it is more important to look at the estimates  $\hat{\theta}_1\hat{\sigma}$  and  $\hat{\theta}_2\hat{\sigma}$  rather than  $\hat{\theta}_1$  and  $\hat{\theta}_2$ . If we look at the estimates  $\hat{\theta}_1\hat{\sigma}$  and  $\hat{\theta}_2\hat{\sigma}$ , this bias pattern disappears.

By comparing the MSEs obtained from the proposed method (I) and the oracle method (III), one can see that, in terms of prediction power, the proposed method is not too far behind from the oracle method. As expected, as sample size increases, the MSE values of the proposed method decrease, which supports the claim that the spline approximation for the fixed effects components  $f_p$ 's is consistent. We also note that the MSE values decrease as SNR increases, which indicates that the noise level affects the convergence rate of the spline approximation. In general, the random effects selection based on BIC improves

TABLE 6  
Simulation results for Model 3. The true numbers of non-zero  $f_p$  and random trends are both 2. The notation  $N_I$  denotes “number of times the random intercept is selected”

SNR = 3		number of non-zero $\hat{f}_p$			MSE			number of random trends			number of times true model is	
configuration	method	2	3	4+	mean	S.E.	$N_I$	$\leq 1$	2	$\geq 3$	included	selected
$n = 128$	I	200	0	0	0.2676	0.0079	147	66	125	9	128	100
$m = 7$	II	106	54	40	0.4079	0.0104	—	—	—	—	—	—
$N = 16$	III	200	0	0	0.2593	0.0086	—	—	—	—	—	—
$n = 256$	I	200	0	0	0.1825	0.0055	190	15	176	9	185	161
$m = 8$	II	150	47	3	0.3093	0.0086	—	—	—	—	—	—
$N = 16$	III	200	0	0	0.1773	0.0056	—	—	—	—	—	—
$n = 512$	I	200	0	0	0.2092	0.0068	199	1	191	8	199	179
$m = 9$	II	173	24	3	0.2989	0.0076	—	—	—	—	—	—
$N = 16$	III	200	0	0	0.1345	0.0061	—	—	—	—	—	—
$n = 256$	I	200	0	0	0.1732	0.0043	134	74	117	9	126	108
$m = 8$	II	155	38	7	0.251	0.0054	—	—	—	—	—	—
$N = 32$	III	200	0	0	0.153	0.0039	—	—	—	—	—	—
$n = 512$	I	200	0	0	0.1254	0.0042	196	4	186	10	196	180
$m = 9$	II	169	29	2	0.1865	0.0054	—	—	—	—	—	—
$N = 32$	III	200	0	0	0.0905	0.0032	—	—	—	—	—	—

SNR = 4		number of non-zero $\hat{f}_p$			MSE			number of random trends			number of times true model is	
configuration	method	2	3	4+	mean	S.E.	$N_I$	$\leq 1$	2	$\geq 3$	included	selected
$n = 128$	I	200	0	0	0.2049	0.0051	153	77	110	13	120	90
$m = 7$	II	120	51	29	0.3013	0.0067	—	—	—	—	—	—
$N = 16$	III	200	0	0	0.1871	0.0047	—	—	—	—	—	—
$n = 256$	I	200	0	0	0.144	0.0037	188	19	171	10	181	154
$m = 8$	II	151	44	5	0.2227	0.0054	—	—	—	—	—	—
$N = 16$	III	200	0	0	0.1274	0.0033	—	—	—	—	—	—
$n = 512$	I	200	0	0	0.1369	0.0045	200	0	194	6	200	180
$m = 9$	II	149	46	5	0.1917	0.0054	—	—	—	—	—	—
$N = 16$	III	200	0	0	0.088	0.0038	—	—	—	—	—	—
$n = 256$	I	200	0	0	0.1338	0.0026	135	77	116	7	123	108
$m = 8$	II	151	40	9	0.1838	0.0038	—	—	—	—	—	—
$N = 32$	III	200	0	0	0.117	0.0022	—	—	—	—	—	—
$n = 512$	I	200	0	0	0.0827	0.0021	198	4	183	13	196	175
$m = 9$	II	157	39	4	0.1183	0.003	—	—	—	—	—	—
$N = 32$	III	200	0	0	0.0646	0.0016	—	—	—	—	—	—

when sample size and/or SNR increase. Lastly, for the repeated measurement model, when the number of replicates per subject increases, our method results in a more accurate random effects selection. However, the MSE values for the fixed effects components do not decrease when the number of replicates per subject increases, unless the total sample size increases.

In conclusion, our proposed method using BIC model selection is very capable of recovering the true model, especially when the sample size and/or SNR are not too small.

TABLE 7  
Estimates of random components for Model 3

SNR = 3		$\hat{\theta}_1 (\theta_1 = 2)$			$\hat{\theta}_2 (\theta_2 = 3)$			$\hat{\theta}_3 (\theta_3 = 4)$			$\hat{\sigma} (\sigma \approx 1)$		
configu- ration	method	median	mean	s.e.	median	mean	s.e.	median	mean	s.e.	median	mean	s.e.
$n = 128$	I	1.979	2.032	0.028	2.857	2.903	0.046	3.412	3.541	0.054	0.849	0.851	0.007
$m = 7$	II	1.83	1.86	0.047	2.899	2.947	0.065	3.568	3.673	0.063	0.821	0.811	0.008
$N = 16$	III	1.832	1.851	0.037	2.837	2.85	0.053	3.561	3.665	0.057	0.798	0.799	0.006
$n = 256$	I	1.891	1.927	0.027	2.806	2.878	0.045	3.666	3.658	0.058	0.829	0.829	0.004
$m = 8$	II	1.969	2.004	0.032	3.026	3.055	0.049	3.612	3.646	0.06	0.839	0.834	0.004
$N = 16$	III	1.86	1.866	0.029	2.769	2.818	0.046	3.699	3.705	0.059	0.81	0.809	0.003
$n = 512$	I	2.026	2.062	0.03	3.074	3.077	0.05	3.779	3.817	0.052	0.806	0.806	0.002
$m = 9$	II	2.136	2.167	0.032	3.234	3.248	0.05	3.819	3.827	0.052	0.807	0.809	0.003
$N = 16$	III	1.866	1.869	0.026	2.82	2.789	0.043	3.78	3.803	0.052	0.797	0.797	0.002
$n = 256$	I	2.041	2.058	0.019	2.767	2.831	0.032	3.728	3.694	0.045	0.854	0.879	0.008
$m = 8$	II	1.851	1.85	0.028	2.792	2.825	0.039	3.778	3.74	0.048	0.829	0.835	0.005
$N = 32$	III	1.908	1.912	0.027	2.794	2.837	0.036	3.881	3.88	0.047	0.801	0.806	0.004
$n = 512$	I	2.007	2.041	0.021	2.922	2.949	0.03	3.787	3.837	0.037	0.808	0.812	0.003
$m = 9$	II	2.055	2.077	0.022	3.017	3.048	0.031	3.786	3.829	0.038	0.806	0.809	0.003
$N = 32$	III	1.974	1.99	0.021	2.861	2.89	0.029	3.832	3.875	0.038	0.797	0.8	0.003

SNR = 4		$\hat{\theta}_1 (\theta_1 = 2)$			$\hat{\theta}_2 (\theta_2 = 3)$			$\hat{\theta}_3 (\theta_3 = 4)$			$\hat{\sigma} (\sigma \approx 0.8)$		
configu- ration	method	median	mean	s.e.	median	mean	s.e.	median	mean	s.e.	median	mean	s.e.
$n = 128$	I	1.839	1.938	0.031	2.689	2.745	0.043	3.363	3.459	0.058	0.661	0.669	0.006
$m = 7$	II	1.709	1.706	0.05	2.564	2.595	0.06	3.367	3.441	0.073	0.659	0.655	0.007
$N = 16$	III	1.773	1.795	0.038	2.616	2.606	0.052	3.532	3.601	0.062	0.621	0.622	0.004
$n = 256$	I	1.815	1.83	0.026	2.712	2.766	0.044	3.545	3.545	0.049	0.638	0.644	0.003
$m = 8$	II	1.839	1.865	0.031	2.88	2.895	0.048	3.506	3.534	0.05	0.648	0.649	0.003
$N = 16$	III	1.784	1.781	0.028	2.69	2.71	0.045	3.618	3.63	0.049	0.623	0.625	0.002
$n = 512$	I	2.039	2.064	0.029	3.09	3.086	0.045	3.754	3.731	0.049	0.615	0.613	0.002
$m = 9$	II	2.166	2.195	0.03	3.303	3.274	0.046	3.796	3.749	0.049	0.615	0.615	0.002
$N = 16$	III	1.871	1.874	0.027	2.796	2.784	0.04	3.75	3.721	0.05	0.609	0.607	0.002
$n = 256$	I	1.969	1.991	0.018	2.746	2.768	0.031	3.446	3.503	0.039	0.664	0.682	0.006
$m = 8$	II	1.825	1.831	0.026	2.748	2.76	0.038	3.471	3.57	0.041	0.642	0.645	0.004
$N = 32$	III	1.901	1.888	0.024	2.78	2.783	0.036	3.647	3.712	0.039	0.62	0.622	0.003
$n = 512$	I	1.926	1.946	0.019	2.888	2.917	0.03	3.653	3.754	0.037	0.616	0.617	0.002
$m = 9$	II	1.977	1.992	0.02	2.995	2.999	0.032	3.665	3.745	0.037	0.616	0.616	0.002
$N = 32$	III	1.9	1.902	0.018	2.806	2.848	0.03	3.721	3.796	0.037	0.609	0.609	0.002

## 7. Real data sets

### 7.1. Prostate cancer data

This data set has been studied by Tibshirani (1996). It examines the correlation between the level of prostate specific antigen and a number of clinical measures in men who were about to receive a radical prostatectomy. Table 8 lists the variables in this data set. The full model is

$$\begin{aligned} \text{lpsa} = & \mu + f_1(\text{lcavol}) + f_2(\text{lweight}) + f_3(\text{age}) + f_4(\text{lbph}) + f_5(\text{gleason}) \\ & + f_6(\text{pgg45}) + I(\text{svi} = 0)u_0 + I(\text{svi} = 1)u_1 + \epsilon, \end{aligned}$$

where  $\mu$  is the overall mean,  $f_p$  are the additive components,  $u_0, u_1 \sim N(0, \theta^2\sigma^2)$  are the random effects corresponding to variable  $\text{svi}$ , and  $\epsilon$  is  $N(0, \sigma^2)$  random error. In this model, we mapped the domains of all continuous predictors into  $[0, 1]$ .

TABLE 8  
Variables of the Prostate Cancer Data Set

Variable Name	Description	Type
<code>lpsa</code>	log of prostate specific antigen	Continuous
<code>lcavol</code>	log cancer volume	Continuous
<code>lweight</code>	log prostate weight	Continuous
<code>age</code>	age	Continuous
<code>lbph</code>	log of benign prostatic hyperplasia amount	Continuous
<code>lcp</code>	log of capsular penetration	Continuous
<code>gleason</code>	gleason score	Continuous
<code>pgg45</code>	Gleason score 4 or 5	Continuous
<code>svi</code>	seminal vesicle invasion	Binary

The analysis by Tibshirani (1996) concludes that variables `lcavol`, `lweight` and `svi` are important in predicting the level of prostate specific antigen. Our proposed method gives similar results: we confirm that `lcavol` and `lweight` are related to `lpsa`, the binary predictor `svi` is significant, and all the remaining variables are insignificant. Plots of the estimated fixed effects  $f_p$ 's and fitted values can be found in Figure 1.

## 7.2. Depression data

It is of interest to examine the degree to which the effects of time-varying drug plasma levels on the change in depression levels over time. In this study each patient was treated for 4 weeks with imipramine. Blood samples were drawn twice per week, 15 hours after the last drug intake, and imipramine and desipramine concentrations were measured. Using this data set, Reisby *et al.* (1977) examine the relationship amongst imipramine (`imi`), desipramine (`dmi`) and clinical response in 66 depressed inpatients (37 endogenous and 29 non-endogenous). The variables of this data set is summarized in Table 9.

The full model is

$$\begin{aligned} \text{hamd}_{ik} = & \mu + u_i + f_1(\text{week}_{ik}) + f_2(\text{week}_{ik} \times \text{endog}_i) + f_3(\text{imi}_{ik}) \\ & + f_4(\text{imi}_{ik} \times \text{endog}_i) + f_5(\text{dmi}_{ik}) + f_6(\text{dmi}_{ik} \times \text{endog}_i) + \epsilon_{ik}, \end{aligned}$$

where  $i = 1, \dots, 66$ ,  $k = 1, \dots, 4$ ,  $\mu$  is the overall mean,  $f_p$ 's are the additive components,  $u_i \sim N(0, \theta^2 \sigma^2)$ 's are the random effects corresponding to different inpatients and  $\epsilon_{ik}$ 's are i.i.d.  $N(0, \sigma^2)$  errors. Note that we introduced a random intercept  $u_i$  for each subject, and interaction terms of `endog` with other continuous predictors. We have also mapped the domains of all the continuous predictors into  $[0, 1]$ .

With this data set, we demonstrate that our proposed method can handle repeated measurement data. The variation among different patients can be explained by a random intercept term. Variable `id` represents an important characteristics of repeated measurement data and is kept in our model. And the group lasso procedure can help to choose the additive components.

Hedeker and Gibbons (2006) have also investigated this data set using linear mixed model and we obtained similar conclusions: larger `dmi` values lead to

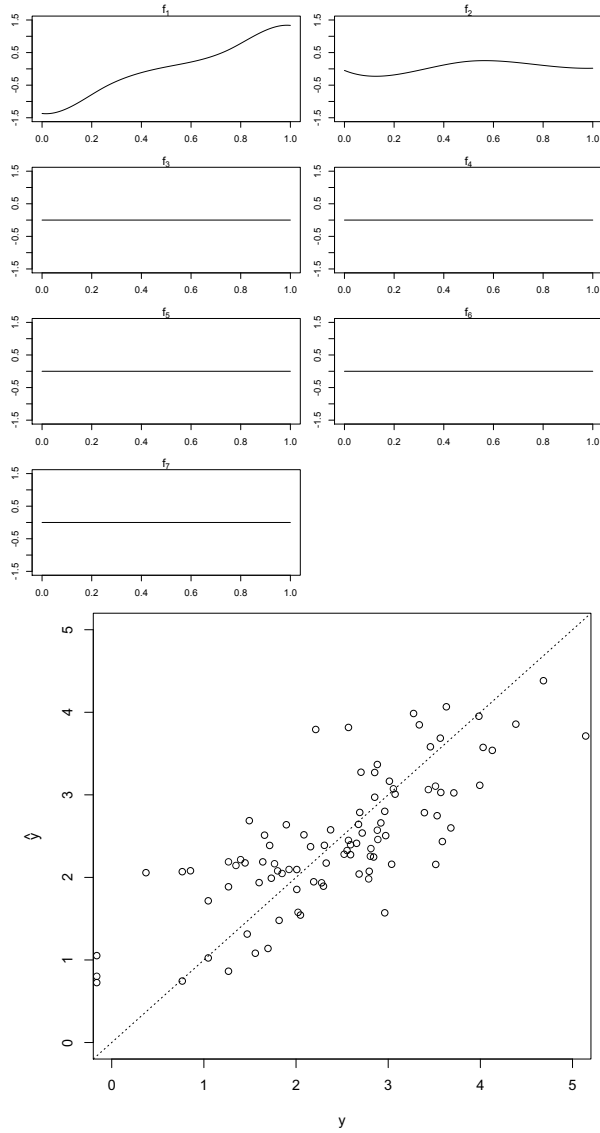


FIG 1. Top: estimated fixed effects  $f_p$ 's of the Prostate Cancer Data Set. Bottom: fitted values versus response.

TABLE 9  
Variables of the Depression Data Set

Variable Name	Description	Type
id	subject number	Discrete (1 to 66)
hamd	Hamilton Depression Scores	Continuous
endog	endogenous (=1) or non-endogenous (=0)	Binary
week	week	Continuous
imi	imipramine drug-plasma levels ( $\mu g/l$ )	Continuous
dmi	desipramine drug-plasma levels ( $\mu g/l$ )	Continuous

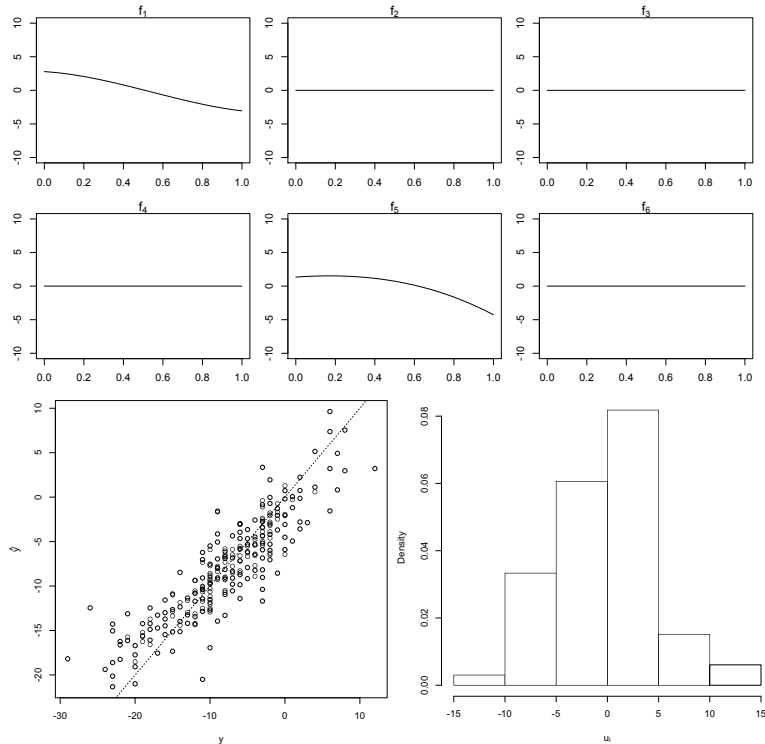


FIG 2. Top two rows: estimated functions of the depression data set. Bottom left: fitted values versus response. Bottom right: histogram for the predicted  $u_i$ .

greater improvement in depression (i.e., more negative `hamd` scores), while the drug `imi` is not significantly related to the `hamd` scores. Plots of estimated fixed effects  $f_p$ 's, histograms of the estimated  $u_i$ 's, together with fitted values can be found in Figure 2.

## 8. Summary

In this paper we studied the problem of variable selection in the context of non-parametric additive mixed modeling. The mixed modeling framework provides a method to jointly handle additive fixed effects as well as random effects. The additive fixed effects component is approximated by truncated series expansions with B-spline bases, with consistency properties of the spline approximation established. We have considered fitting the nonparametric fixed components with the adaptive group lasso methodology. This adaptive group lasso procedure yields sparse estimates for the coefficients corresponding to the B-spline bases. An model selection criterion derived from BIC, which is proven to be consistent, is presented for selecting the tuning parameters ( $\lambda$  and  $\nu$ ) of the adaptive group lasso procedure. Empirical results show that the new methodology is very

efficient in fitting high dimensional data with sparse solution, especially for the case of repeated measurement model with a number of continuous covariates.

### Acknowledgment

The authors are grateful to the reviewer and the editor for many constructive comments, which led to a much improved version of the paper.

### Appendix

#### Appendix A: Derivation of (11) and (12)

This appendix derives (11) and (12). In below  $\otimes$  denotes the Kronecker product operator and  $\text{vec}\mathbf{X}$  denotes the vector obtained by stacking the columns of the matrix  $\mathbf{X}$ . We begin by calculating

$$\begin{aligned}\frac{\partial \mathbf{r}'\mathbf{V}^{*-1}\mathbf{r}}{\partial \text{vec}\mathbf{G}} &= \frac{\partial \text{vec}\mathbf{V}^*}{\partial \text{vec}\mathbf{G}} \frac{\partial \mathbf{r}'\mathbf{V}^{*-1}\mathbf{r}}{\partial \text{vec}\mathbf{V}^*} \\ &= -\mathbf{Z}'\mathbf{V}^{*-1}\mathbf{r} \otimes \mathbf{Z}'\mathbf{V}^{*-1'}\mathbf{r} \\ &= -\text{vec}(\mathbf{Z}'\mathbf{V}^{*-1'}\mathbf{r}\mathbf{r}'\mathbf{V}^{*-1'}\mathbf{Z}),\end{aligned}$$

$$\begin{aligned}\frac{\partial^2 \mathbf{r}'\mathbf{V}^{*-1}\mathbf{r}}{\partial \text{vec}(\mathbf{G}')\partial \text{vec}(\mathbf{G})'} &= (\mathbf{Z}' \otimes \mathbf{Z}') \left\{ (\mathbf{V}^{*-1'}\mathbf{r} \otimes \mathbf{V}^{*-1}\mathbf{Z})(\mathbf{r}'\mathbf{V}^{*-1'}\mathbf{Z} \otimes \mathbf{I}) \right. \\ &\quad \left. + (\mathbf{V}^{*-1'}\mathbf{Z} \otimes \mathbf{V}^{*-1}\mathbf{r})(\mathbf{I} \otimes \mathbf{r}'\mathbf{V}^{*-1}\mathbf{Z}) \right\} \\ &= \mathbf{Z}'\mathbf{V}^{*-1'}\mathbf{r}\mathbf{r}'\mathbf{V}^{*-1'}\mathbf{Z} \otimes \mathbf{Z}'\mathbf{V}^{*-1}\mathbf{Z} \\ &\quad + \mathbf{Z}'\mathbf{V}^{*-1'}\mathbf{Z} \otimes \mathbf{Z}'\mathbf{V}^{*-1}\mathbf{r}\mathbf{r}'\mathbf{V}^{*-1}\mathbf{Z},\end{aligned}$$

$$\frac{\partial \log(|\mathbf{V}^*|)}{\partial \text{vec}\mathbf{G}} = (\mathbf{Z}' \otimes \mathbf{Z}')\text{vec}(\mathbf{V}^{*-1'}) = \text{vec}(\mathbf{Z}'\mathbf{V}^{*-1'}\mathbf{Z})$$

and

$$\begin{aligned}\frac{\partial^2 \log(|\mathbf{V}^*|)}{\partial \text{vec}(\mathbf{G}')\partial \text{vec}(\mathbf{G})'} &= -(\mathbf{Z}' \otimes \mathbf{Z}')(\mathbf{V}^{*-1'}\mathbf{Z} \otimes \mathbf{V}^{*-1}\mathbf{Z}) \\ &= -\mathbf{Z}'\mathbf{V}^{*-1'}\mathbf{Z} \otimes \mathbf{Z}'\mathbf{V}^{*-1}\mathbf{Z}.\end{aligned}$$

Let  $\mathbf{A} = \mathbf{Z}'\mathbf{V}^{*-1}\mathbf{r}(\mathbf{Z}'\mathbf{V}^{*-1}\mathbf{r})'$  and  $\mathbf{B} = \mathbf{Z}'\mathbf{V}^{*-1}\mathbf{Z}$ . Since  $\mathbf{G}$  is symmetric, the gradient and Hessian matrix w.r.t.  $\text{vec}\mathbf{G}$  are

$$\begin{aligned}\frac{\partial p(\mathbf{G}|\hat{\beta})}{\partial \text{vec}(\mathbf{G})} &= \text{vec}(\mathbf{Z}'\mathbf{V}^{*-1'}\mathbf{Z}) - \frac{n}{\mathbf{r}'\mathbf{V}^{*-1}\mathbf{r}}\text{vec}(\mathbf{Z}'\mathbf{V}^{*-1'}\mathbf{r}\mathbf{r}'\mathbf{V}^{*-1'}\mathbf{Z}) \\ &= \text{vec}\mathbf{B} - \frac{1}{\hat{\sigma}^2}\text{vec}\mathbf{A}\end{aligned}$$

and

$$\frac{\partial^2 p(\mathbf{G}|\hat{\boldsymbol{\beta}})}{\partial \text{vec}(\mathbf{G}) \partial \text{vec}(\mathbf{G})'} = \frac{\partial \text{vec}(\mathbf{G}')}{\partial \text{vec}(\mathbf{G})} \frac{\partial^2 p(\mathbf{G}|\hat{\boldsymbol{\beta}})}{\partial \text{vec}(\mathbf{G}') \partial \text{vec}(\mathbf{G})'},$$

where

$$\begin{aligned} & \frac{\partial^2 p(\mathbf{G}|\hat{\boldsymbol{\beta}})}{\partial \text{vec}(\mathbf{G}') \partial \text{vec}(\mathbf{G})'} \\ &= -\mathbf{Z}' \mathbf{V}^{*-1'} \mathbf{Z} \otimes \mathbf{Z}' \mathbf{V}^{*-1} \mathbf{Z} \\ & \quad - \frac{n}{(\mathbf{r}' \mathbf{V}^{*-1} \mathbf{r})^2} \text{vec}(\mathbf{Z}' \mathbf{V}^{*-1'} \mathbf{r} \mathbf{r}' \mathbf{V}^{*-1'} \mathbf{Z}) \text{vec}(\mathbf{Z}' \mathbf{V}^{*-1'} \mathbf{r} \mathbf{r}' \mathbf{V}^{*-1'} \mathbf{Z})' \\ & \quad + \frac{n}{\mathbf{r}' \mathbf{V}^{*-1} \mathbf{r}} \left( \mathbf{Z}' \mathbf{V}^{*-1'} \mathbf{r} \mathbf{r}' \mathbf{V}^{*-1'} \mathbf{Z} \otimes \mathbf{Z}' \mathbf{V}^{*-1} \mathbf{Z} \right. \\ & \quad \left. + \mathbf{Z}' \mathbf{V}^{*-1'} \mathbf{Z} \otimes \mathbf{Z}' \mathbf{V}^{*-1} \mathbf{r} \mathbf{r}' \mathbf{V}^{*-1} \mathbf{Z} \right) \\ &= -\mathbf{B} \otimes \mathbf{B} - \frac{1}{n \hat{\sigma}^4} \text{vec}(\mathbf{A}) \text{vec}(\mathbf{A})' + \frac{1}{\hat{\sigma}^2} (\mathbf{A} \otimes \mathbf{B} + \mathbf{B} \otimes \mathbf{A}). \end{aligned}$$

With the above (11) and (12) can be obtained.

## Appendix B: Lemmas

This appendix presents a few lemmas which will be required for proofing our major theoretical results.

**Lemma 1.** Under (A.3) and (A.4), for every  $f_p \in \mathcal{F}$ , there exists  $\tilde{f}_p \in \mathcal{S}(d, \mathbf{t}_p)$  s.t.

$$\frac{1}{n} \sum_{i=1}^n \left\{ \tilde{f}_p(X_{ip}) - f_p(X_{ip}) \right\}^2 = O_p(m^{-d}).$$

*Proof.* See lemma 1 of Huang *et al.* (2010).  $\square$

**Lemma 2.** Let  $\mathbf{w} = (w_p \frac{\hat{\boldsymbol{\beta}}_p'}{\|\hat{\boldsymbol{\beta}}_p\|}, p \in A_T)'$ . Under (A.1) and (A.5)(b) we have  $\|\mathbf{w}\| = o_p(1)$ .

*Proof.*

$$\|\mathbf{w}\|^2 = \sum_{p \in A_T} w_p^2 = \sum_{p \in A_T} (\|\tilde{\boldsymbol{\beta}}_p\|^{-1} - \nu)_+^2 \leq |A_T| \max_{p \in A_T} (\|\tilde{\boldsymbol{\beta}}_p\|^{-1} - \nu)_+^2 = o_p(1).$$

$\square$

**Lemma 3.** For  $p = 1, \dots, P, j = 1, \dots, m$ , let  $\mathbf{a}_{pj} = [\phi_{pj}(X_{ip})]_{i=1}^n$  and  $\mathbf{e} \sim N(\mathbf{0}, \mathbf{V})$ , then

$$\|\mathbf{a}_{pj}\| = O_p\left(\sqrt{\frac{n}{m}}\right), \quad \mathbf{a}_{pj}' \mathbf{V}^{-1} \mathbf{e} = O_p\left(\sqrt{\frac{n}{m}}\right) \quad \text{and} \quad \|\mathbf{X}_p' \mathbf{V}^{-1} \mathbf{e}\| = O_p(\sqrt{n}).$$



*Proof.* By properties of B-splines, there exist  $c_1$  and  $c_2$  s.t.

$$\|\mathbf{a}_{pj}\|_\infty \leq \|[\phi_{pj}(X_{ip})]_{i=1}^n\|_\infty \leq c_1 \quad \text{and} \quad \mathbb{E}[\phi_{pj}^2(X_{ip})] \leq \frac{c_2}{m}.$$

Then,

$$\phi_{pj}^2(X_{ip}) - \mathbb{E}\phi_{pj}^2(X_{ip}) \leq 2c_1^2 \quad \text{and} \quad \text{Var}[\phi_{pj}^2(X_{ip})] \leq 4c_1^4 \frac{c_2}{m}. \quad (20)$$

By law of large number,

$$\frac{1}{n} \|\mathbf{a}_{pj}\|^2 = \frac{1}{n} \sum_{i=1}^n \phi_{pj}^2(X_{ip}) = O_p\left(\frac{1}{m}\right).$$

Therefore,

$$\mathbf{a}'_{pj} \mathbf{V}^{-1} \mathbf{e} = O_p\left(\sqrt{\frac{n}{m}}\right) \quad \text{and} \quad \|\mathbf{X}'_p \mathbf{V}^{-1} \mathbf{e}\| \leq \sqrt{m} \max_{1 \leq j \leq m} \mathbf{a}'_{pj} \mathbf{V}^{-1} \mathbf{e} = O_p(\sqrt{n}).$$

The proof completes.  $\square$

**Lemma 4.** Let  $\mathbf{X}_A$  be the design matrix corresponding to  $A$ . If  $|A|$  is bounded, then

$$\|\mathbf{X}_A\| = O_p\left(\sqrt{\frac{n}{m}}\right), \quad \text{and} \quad \|(\mathbf{X}'_A \mathbf{V}^{-1} \mathbf{X}_A)^{-1}\| = O_p\left(\frac{m}{n}\right).$$

*Proof.* The result holds under (A.2) and (A.3). Please referring lemma 3 of Huang *et al.* (2010) and lemma 6.2 of Zhou *et al.* (1998) for details.  $\square$

**Lemma 5.** For  $p = 1, \dots, P, j = 1, \dots, m$ , let  $\mathbf{a}_{pj} = [\phi_{pj}(X_{ip})]_{i=1}^n$ . Then

$$\mathbb{E}\left[\max_{pj} \|\mathbf{a}_{pj}\|\right] = O_p\left(\sqrt{\frac{n}{m}}\right).$$

*Proof.* By equation (20) and lemma A.1 of Van De Geer (2008),

$$\mathbb{E}\left[\max_{pj} \left|\frac{1}{n} \sum_{i=1}^n \{\phi_{pj}^2(X_{ip}) - \mathbb{E}\phi_{pj}^2(X_{ip})\}\right|\right] \leq \sqrt{\frac{8c_1^4 c_2 \log(2mP)}{mn}} + 2c_1^2 \frac{\log(2mP)}{n}.$$

Also,

$$\begin{aligned} \mathbb{E}\left[\frac{1}{n} \max_{pj} \|\mathbf{a}_{pj}\|^2\right] &= \mathbb{E}\left[\max_{pj} \left|\frac{1}{n} \sum_{i=1}^n \{\phi_{pj}^2(X_{ip})\}\right|\right] \\ &\leq \mathbb{E}\left[\max_{pj} \left|\frac{1}{n} \sum_{i=1}^n \{\phi_{pj}^2(X_{ip}) - \mathbb{E}\phi_{pj}^2(X_{ip})\}\right|\right] + \max_{pj} \mathbb{E}\phi_{pj}^2(X_{ip}) \end{aligned}$$

and

$$\begin{aligned} \mathbb{E}\left[\max_{pj} \|\mathbf{a}_{pj}\|\right] &\leq \sqrt{\mathbb{E}\left[\max_{pj} \|\mathbf{a}_{pj}\|^2\right]} \\ &\leq \left\{ \sqrt{\frac{8c_1^4 c_2 n \log(2mP)}{m}} + 2c_1^2 \log(2mP) + \frac{c_2 n}{m} \right\}^{1/2}. \end{aligned}$$

Particularly, under (A.6)(a), we have  $\frac{m \log(2mP)}{n} \rightarrow 0$ , so  $E[\max_{pj} \|\mathbf{a}_{pj}\|] = O_p(\sqrt{\frac{n}{m}})$ . Therefore, the result follows.  $\square$

**Lemma 6.** Let  $\mathbf{e} \sim N(\mathbf{0}, \mathbf{V})$ , for any  $A \subset \{1, \dots, P\}$

$$E(\max_{p \in A} \|\mathbf{X}'_p \mathbf{V}^{-1} \mathbf{e}\|) = O_p(\sqrt{n \log(2mP)})$$

*Proof.* For  $p = 1, \dots, P, j = 1, \dots, m$ , let  $\mathbf{a}_{pj} = [\phi_{pj}(X_{ip})]_{i=1}^n$ . Condition on  $\mathbf{a}_{pj}$ , by lemma 2.2.1 and 2.2.2 of Van der Vaart and Wellner (1996),

$$E(\max_{1 \leq p \leq P, 1 \leq j \leq m} \mathbf{a}'_{pj} \mathbf{V}^{-1} \mathbf{e} | \mathbf{a}_{pj}) \leq C_1 \sqrt{\log(2mP)} \cdot \max_{pj} \|\mathbf{a}_{pj} \mathbf{V}^{-1/2}\|,$$

for some  $C_1 > 0$ . Therefore, for some  $C_2 > 0$ ,

$$\begin{aligned} E(\max_{1 \leq p \leq P, 1 \leq j \leq m} \mathbf{a}'_{pj} \mathbf{V}^{-1} \mathbf{e}) &\leq C_1 \sqrt{\log(2mP)} \cdot E\left[\max_{pj} \|\mathbf{a}_{pj}\|\right] \cdot \|\mathbf{V}^{-1/2}\| \\ &\leq C_2 \sqrt{\frac{n \log(2mP)}{m}}. \end{aligned}$$

Moreover,

$$E(\max_{p \in A} \|\mathbf{X}'_p \mathbf{V}^{-1} \mathbf{e}\|) \leq \sqrt{m} \cdot E(\max_{1 \leq p \leq P, 1 \leq j \leq m} \mathbf{a}'_{pj} \mathbf{V}^{-1} \mathbf{e}) \leq C_2 \sqrt{n \log(2mP)}.$$

It completes the proof.  $\square$

### Appendix C: Proofs of theorems

*Proof of theorem 1.* The proof follows that main arguments of Theorem 3 of Huang *et al.* (2010), but with some modified technical details. Recall the necessary and sufficient conditions for  $\hat{\boldsymbol{\beta}} = (\hat{\boldsymbol{\beta}}'_1, \dots, \hat{\boldsymbol{\beta}}'_P)'$  to be the solution of (5) are

$$\begin{cases} -\mathbf{X}'_p \mathbf{V}^{-1}(\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}) + \lambda w_p \frac{\hat{\boldsymbol{\beta}}_p}{\|\hat{\boldsymbol{\beta}}_p\|} = \mathbf{0}, & \text{if } \hat{\boldsymbol{\beta}}_p \neq \mathbf{0}; \\ \|\mathbf{X}'_p \mathbf{V}^{-1}(\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}})\| \leq \lambda w_p, & \text{if } \hat{\boldsymbol{\beta}}_p = \mathbf{0}. \end{cases} \quad (21)$$

Let  $\mathbf{w} = (w_p \frac{\hat{\boldsymbol{\beta}}'_p}{\|\hat{\boldsymbol{\beta}}_p\|}, p \in A_T)'$ . Define

$$\hat{\boldsymbol{\gamma}}_{A_T} = (\hat{\boldsymbol{\gamma}}'_{A_T, p}, p \in A_T)' = (\mathbf{X}'_{A_T} \mathbf{V}^{-1} \mathbf{X}_{A_T})^{-1} (\mathbf{X}'_{A_T} \mathbf{V}^{-1} \mathbf{Y} - \lambda \mathbf{w}).$$

If

$$\begin{cases} \|\hat{\boldsymbol{\gamma}}_{A_T, p}\| > 0, & \forall p \in A_T; \\ \|\mathbf{X}'_p \mathbf{V}^{-1}(\mathbf{Y} - \mathbf{X}_{A_T} \hat{\boldsymbol{\gamma}}_{A_T})\| \leq \lambda w_p, & \forall p \in A_0, \end{cases}$$

then  $\hat{\boldsymbol{\beta}} = (\hat{\boldsymbol{\gamma}}'_{A_T}, \mathbf{0}')'$  satisfies (21),  $\mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{X}_{A_T} \hat{\boldsymbol{\gamma}}_{A_T}$  and  $\text{sign}(\|\hat{\boldsymbol{\beta}}_p\|) = \text{sign}(\|\boldsymbol{\beta}_p\|)$  for  $p = 1, \dots, P$ . Therefore,  $\text{sign}(\|\hat{\boldsymbol{\beta}}_p\|) = \text{sign}(\|\boldsymbol{\beta}_p\|)$  for  $p = 1, \dots, P$  if

$$\begin{cases} \|\boldsymbol{\beta}_p\| - \|\hat{\boldsymbol{\gamma}}_{A_T, p}\| < \|\boldsymbol{\beta}_p\|, & \forall p \in A_T; \\ \|\mathbf{X}'_p \mathbf{V}^{-1}(\mathbf{Y} - \mathbf{X}_{A_T} \hat{\boldsymbol{\gamma}}_{A_T})\| \leq \lambda w_p, & \forall p \in A_0. \end{cases}$$

Now

$$\begin{aligned} & P(\text{sign}(\|\hat{\boldsymbol{\beta}}_p\|) \neq \text{sign}(\|\boldsymbol{\beta}_p\|), \exists p) \\ & \leq P(\|\hat{\boldsymbol{\gamma}}_{A_T,p} - \boldsymbol{\beta}_p\| \geq \|\boldsymbol{\beta}_p\|, \exists p \in A_T) \\ & \quad + P(\|\mathbf{X}'_p \mathbf{V}^{-1}(\mathbf{Y} - \mathbf{X}_{A_T} \hat{\boldsymbol{\gamma}}_{A_T})\| > \lambda w_p, \exists p \in A_0). \end{aligned} \quad (22)$$

Next we show that the right hand side of (22) tends to zero, by using claim 1 and claim 2 below.

**Claim 1:**

$$P(\|\hat{\boldsymbol{\gamma}}_{A_T,p} - \boldsymbol{\beta}_p\| \geq \|\boldsymbol{\beta}_p\|, \exists p \in A_T) \rightarrow 0. \quad (23)$$

$$\text{Let } \mathbf{T}_p = [\mathbf{T}_{p1}, \mathbf{T}_{p2}, \dots, \mathbf{T}_{p|A_T|}], \text{ where } \mathbf{T}_{pj} = \begin{cases} \mathbf{I}, & p = j; \\ \mathbf{0}, & p \neq j. \end{cases}$$

And let  $\mathbf{C}_A = \mathbf{X}'_A \mathbf{V}^{-1} \mathbf{X}_A$ ,  $\mathbf{e} = \mathbf{Z}\mathbf{u} + \boldsymbol{\epsilon}$ ,  $\boldsymbol{\delta} = [\delta_i]_{i=1}^n$ , where  $\delta_i = \sum_{p=1}^P f(\mathbf{X}_{ip}) - \mathbf{X}_{(i)}\boldsymbol{\beta}$  and  $\mathbf{X}_{(i)}$  is the  $i^{\text{th}}$  row of  $\mathbf{X}$ . We have

$$\hat{\boldsymbol{\gamma}}_{A_T,p} - \boldsymbol{\beta}_p = \mathbf{T}_p \mathbf{C}_{A_T}^{-1} \left\{ \mathbf{X}'_{A_T} \mathbf{V}^{-1}(\mathbf{e} + \boldsymbol{\delta}) - \lambda \mathbf{w} \right\}.$$

By the triangle inequality,

$$\|\hat{\boldsymbol{\gamma}}_{A_T,p} - \boldsymbol{\beta}_p\| \leq \|\mathbf{T}_p \mathbf{C}_{A_T}^{-1} \mathbf{X}'_{A_T} \mathbf{V}^{-1} \mathbf{e}\| + \|\mathbf{T}_p \mathbf{C}_{A_T}^{-1} \mathbf{X}'_{A_T} \mathbf{V}^{-1} \boldsymbol{\delta}\| + \lambda \|\mathbf{T}_p \mathbf{C}_{A_T}^{-1} \mathbf{w}\|.$$

By lemma 4,  $\|\mathbf{C}_{A_T}^{-1}\| = O_p(m/n)$  and  $\|\mathbf{X}_{A_T}\| = O_p(\sqrt{n/m})$ . In addition, by lemma 3, the first term

$$\|\mathbf{T}_p \mathbf{C}_{A_T}^{-1} \mathbf{X}'_{A_T} \mathbf{V}^{-1} \mathbf{e}\| \leq \|\mathbf{C}_{A_T}^{-1}\| \cdot \|\mathbf{X}'_{A_T} \mathbf{V}^{-1} \mathbf{e}\| = O_p(m/\sqrt{n}) = o_p(1).$$

By lemma 1, the second term

$$\|\mathbf{T}_p \mathbf{C}_{A_T}^{-1} \mathbf{X}'_{A_T} \mathbf{V}^{-1} \boldsymbol{\delta}\| \leq \|\mathbf{C}_{A_T}^{-1}\| \cdot \|\mathbf{X}'_{A_T} \mathbf{V}^{-1}\| \cdot \|\boldsymbol{\delta}\| = O_p(m^{-d+1/2}) = o_p(1).$$

By lemma 2, the third term

$$\lambda \|\mathbf{T}_p \mathbf{C}_{A_T}^{-1} \mathbf{w}\| \leq \lambda \|\mathbf{C}_{A_T}^{-1}\| \cdot \|\mathbf{w}\| = O_p(\lambda m/n) o_p(1) = o_p(1).$$

With probability tending to one,  $\|\hat{\boldsymbol{\gamma}}_{A_T,p} - \boldsymbol{\beta}_p\| = o_p(1)$ . Thus, claim 1 is proved.

**Claim 2:**

$$P(\|\mathbf{X}'_p \mathbf{V}^{-1}(\mathbf{Y} - \mathbf{X}_{A_T} \hat{\boldsymbol{\gamma}}_{A_T})\| > \lambda w_p, \exists p \in A_0) \rightarrow 0. \quad (24)$$

Let  $\mathbf{H} = \mathbf{I} - \mathbf{X}_{A_T} \mathbf{C}_{A_T}^{-1} \mathbf{X}'_{A_T} \mathbf{V}^{-1}$ . We have  $\|\mathbf{H}\| = O_p(1)$  and

$$\frac{1}{\lambda w_p} \mathbf{X}'_p \mathbf{V}^{-1}(\mathbf{Y} - \mathbf{X}_{A_T} \hat{\boldsymbol{\gamma}}_{A_T}) = \frac{1}{\lambda w_p} \mathbf{X}'_p \mathbf{V}^{-1}(\mathbf{H}\mathbf{e} + \mathbf{H}\boldsymbol{\delta} + \lambda \mathbf{X}_{A_T} \mathbf{C}_{A_T}^{-1} \mathbf{w}).$$

The second term

$$\frac{1}{\lambda w_p} \|\mathbf{X}'_p \mathbf{V}^{-1} \mathbf{H}\boldsymbol{\delta}\| \leq \frac{1}{\lambda w_p} \|\mathbf{X}'_p \mathbf{V}^{-1}\| \cdot \|\mathbf{H}\| \cdot \|\boldsymbol{\delta}\| = O_p\left(\frac{n}{\lambda r m^{d+1/2}}\right) = o_p(1).$$

The third term

$$\frac{1}{w_p} \|\mathbf{X}'_p \mathbf{V}^{-1} \mathbf{X}_{A_T} \mathbf{C}_{A_T}^{-1} \mathbf{w}\| \leq \frac{1}{w_p} \|\mathbf{X}'_p \mathbf{V}^{-1} \mathbf{X}_{A_T} \mathbf{C}_{A_T}^{-1}\| \|\mathbf{w}\| = O_p(r^{-1}) o_p(1) = o_p(1).$$

By lemma 6,

$$\begin{aligned} & P(\|\mathbf{X}'_p \mathbf{V}^{-1}(\mathbf{Y} - \mathbf{X}_{A_T} \hat{\boldsymbol{\gamma}}_{A_T})\| > \lambda w_p, \exists p \in A_0) \\ & \rightarrow P\left(\frac{1}{\lambda w_p} \|\mathbf{X}'_p \mathbf{V}^{-1} \mathbf{H} \mathbf{e}\| > 1, \exists p \in A_0\right) \\ & = P\left(\max_{p \in A_0} \frac{1}{\lambda w_p} \|\mathbf{X}'_p \mathbf{V}^{-1} \mathbf{H} \mathbf{e}\| > 1\right) \\ & \leq \mathbb{E}\left[\max_{p \in A_0} \frac{1}{\lambda w_p} \|\mathbf{X}'_p \mathbf{V}^{-1} \mathbf{H} \mathbf{e}\|\right] = O_p\left(\frac{\sqrt{n \log(2mP)}}{\lambda r}\right) = o_p(1). \end{aligned}$$

Thus, claim 2 is proved. This together with claim 1 proves that

$$P\left(\text{sign}(\|\hat{\boldsymbol{\beta}}_p\|) = \text{sign}(\|\boldsymbol{\beta}_p\|), p = 1, \dots, P\right) \rightarrow 1.$$

□

*Proof of part 2 of theorem 1.* Define the event

$$F \equiv \bigcap_{p \in A_0} \{\|\mathbf{X}'_p \mathbf{V}^{-1}(\mathbf{Y} - \mathbf{X} \hat{\mathbf{b}}_{A_T})\| \leq \lambda w_p\} \cap \left\{ \min_{p \in A_T} \|\tilde{\boldsymbol{\beta}}_p\| \geq \nu^{-1} \right\}.$$

It is straightforward to check that  $\hat{\mathbf{b}}_{A_T}$  satisfies (21) on  $F$ . Therefore,

$$\begin{aligned} P(\hat{\boldsymbol{\beta}} \neq \hat{\mathbf{b}}_{A_T}) &= P(\|\mathbf{X}'_p \mathbf{V}^{-1}(\mathbf{Y} - \mathbf{X} \hat{\mathbf{b}}_{A_T})\| \\ &> \lambda w_p, \exists p \in A_0) + P\left(\min_{p \in A_T} \|\tilde{\boldsymbol{\beta}}_p\| < \nu^{-1}\right). \end{aligned} \quad (25)$$

By a similar argument as claim 2, put  $\mathbf{w} = \mathbf{0}$ , we can show that the first term goes to zero as  $n \rightarrow \infty$ . This together with (A.5b) implies that the righthand side of (25) goes to zero as  $n \rightarrow \infty$ . This completes the proof of part 2. □

*Proof of parts 3 and 4 of theorem 1.* Since  $\hat{\boldsymbol{\beta}}$  is the minimizer of the quantity  $\frac{1}{2}(\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})' \mathbf{V}^{-1}(\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}) + \lambda \sum_{p=1}^P w_p \|\boldsymbol{\beta}_p\|$ , we have

$$\frac{1}{2} \left\| \mathbf{V}^{-1/2}(\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}) \right\|^2 + \lambda \sum_{p=1}^P w_p \|\hat{\boldsymbol{\beta}}_p\| \leq \frac{1}{2} \left\| \mathbf{V}^{-1/2}(\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}) \right\|^2 + \lambda \sum_{p=1}^P w_p \|\boldsymbol{\beta}_p\|.$$

Therefore,

$$\frac{1}{2} \left\| \mathbf{V}^{-1/2}(\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}) \right\|^2 - \frac{1}{2} \left\| \mathbf{V}^{-1/2}(\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}) \right\|^2 \leq \lambda \sum_{p=1}^P w_p (\|\boldsymbol{\beta}_p\| - \|\hat{\boldsymbol{\beta}}_p\|).$$

Since  $\|\beta_p\| = 0$  for  $p \in A_0$ , the left hand side

$$\begin{aligned} \lambda \sum_{p=1}^P w_p (\|\beta_p\| - \|\hat{\beta}_p\|) &\leq \lambda \sum_{p \in A_T} w_p (\|\beta_p\| - \|\hat{\beta}_p\|) \\ &\leq \lambda \left( \sum_{p \in A_T} w_p^2 \right)^{1/2} \left( \sum_{p \in A_T} \|\beta_p - \hat{\beta}_p\|^2 \right)^{1/2} \\ &\leq \lambda \left( \sum_{p \in A_T} w_p^2 \right)^{1/2} \|\hat{\beta} - \beta\|. \end{aligned}$$

Recall that  $\mathbf{Y} - \mathbf{X}\beta = \mathbf{e} + \delta$ . On the other side, let  $\boldsymbol{\eta} = \mathbf{V}^{-1/2}(\mathbf{e} + \delta)$ ,  $\boldsymbol{\xi} = \mathbf{V}^{-1/2}\mathbf{X}(\hat{\beta} - \beta)$ ,

$$\mathbf{V}^{-1/2}(\mathbf{Y} - \mathbf{X}\hat{\beta}) = \boldsymbol{\eta} - \boldsymbol{\xi}$$

and

$$\|\mathbf{V}^{-1/2}(\mathbf{Y} - \mathbf{X}\hat{\beta})\|^2 = \|\boldsymbol{\eta}\|^2 - 2\boldsymbol{\eta}'\boldsymbol{\xi} + \|\boldsymbol{\xi}\|^2.$$

Therefore,

$$\|\boldsymbol{\xi}\|^2 \leq 2\lambda \left( \sum_{p \in A_T} w_p^2 \right)^{1/2} \|\hat{\beta} - \beta\| + 2\boldsymbol{\eta}'\boldsymbol{\xi}. \tag{26}$$

Let  $A_1 = \{p : \|\beta_p\| > 0 \text{ or } \|\hat{\beta}_p\| > 0\}$ ,  $\mathbf{P} = \mathbf{V}^{-1/2}\mathbf{X}_{A_1}\mathbf{C}_{A_1}^{-1}\mathbf{X}'_{A_1}\mathbf{V}^{-1/2}$ ,  $\mathbf{P}\boldsymbol{\xi} = \boldsymbol{\xi}$ . By Cauchy-Schwartz inequality and completing square,

$$2|\boldsymbol{\eta}'\boldsymbol{\xi}| = 2|\boldsymbol{\eta}'\mathbf{P}\boldsymbol{\xi}| \leq 2\|\mathbf{P}\boldsymbol{\eta}\| \cdot \|\boldsymbol{\xi}\| \leq 2\|\mathbf{P}\boldsymbol{\eta}\|^2 + \frac{1}{2}\|\boldsymbol{\xi}\|^2.$$

Then,

$$\|\boldsymbol{\xi}\|^2 \leq 4\lambda \left( \sum_{p \in A_T} w_p^2 \right)^{1/2} \|\hat{\beta} - \beta\| + 4\|\mathbf{P}\boldsymbol{\eta}\|^2.$$

There exists  $C > 0$ , s.t.  $\|\boldsymbol{\xi}\|^2 \geq C\frac{n}{m}\|\hat{\beta} - \beta\|^2$ . By Completing square,

$$\begin{aligned} \frac{nC}{m}\|\hat{\beta} - \beta\|^2 &\leq \frac{4m}{nC}\lambda^2 \left( \sum_{p \in A_T} w_p^2 \right) + \frac{nC}{2m}\|\hat{\beta} - \beta\|^2 + 4\|\mathbf{P}\boldsymbol{\eta}\|^2 \\ &\leq \frac{8m}{nC}\lambda^2 \left( \sum_{p \in A_T} w_p^2 \right) + 8\|\mathbf{P}\boldsymbol{\eta}\|^2 \\ \|\hat{\beta} - \beta\|^2 &\leq \frac{8m^2}{n^2C^2}\lambda^2 \left( \sum_{p \in A_T} w_p^2 \right) + \frac{8m}{nC}\|\mathbf{P}\boldsymbol{\eta}\|^2. \end{aligned}$$

By lemma 2,  $\sum_{p \in A_T} w_p^2 = o_p(1)$ . Since  $\sqrt{|A_1|} = O_p(1)$ , by similar arguments in claim 1,

$$\|\mathbf{P}\boldsymbol{\eta}\|^2 \leq \|\mathbf{P}\mathbf{V}^{-1/2}\mathbf{e}\|^2 + \|\mathbf{P}\mathbf{V}^{-1/2}\delta\|^2 = O_p(m) + O_p(m^{-2d}).$$

So,

$$\sum_{p=1}^P \|\hat{\beta}_p - \beta_p\|^2 = \|\hat{\beta} - \beta\|^2 = O_p\left(\frac{m^2}{n}\right) + O_p\left(\frac{1}{m^{2d-1}}\right) + O_p\left(\frac{m^2\lambda^2}{n^2}\right) o_p(1).$$

By spline properties, there exist constants  $c_1$  and  $c_2 > 0$ ,

$$c_1 m^{-1} \|\hat{\beta}_p - \beta_p\|^2 \leq \|\hat{f}_p - f\|^2 \leq c_2 m^{-1} \|\hat{\beta}_p - \beta_p\|^2.$$

Therefore,

$$\sum_{p=1}^P \|\hat{f}_p - f_p\|^2 = O_p\left(\frac{m}{n}\right) + O_p\left(\frac{1}{m^{2d}}\right) + O_p\left(\frac{m^2\lambda^2}{n^2}\right) o_p(1).$$

□

*Proof of theorem 2.* We adopt the idea of Wang *et al.* (2007a). Let  $\Omega_- \equiv \{(\lambda, \nu) : \hat{A}_{\lambda, \nu} \cap A_T \neq A_T\}$  and  $\Omega_+ \equiv \{(\lambda, \nu) : \hat{A}_{\lambda, \nu} \supseteq A_T\}$ . It is enough to show that as  $n \rightarrow \infty$ , the following hold:

$$P\left(\text{BIC}(\lambda_n^*, \nu_n^*) = \log(\hat{\sigma}_{A_T}^2) + \frac{1}{n\hat{\sigma}_{A_T}^2} \|\mathbf{Y} - \mathbf{X}\hat{\mathbf{b}}_{A_T} - \mathbf{Z}\hat{\mathbf{u}}_{A_T}\|^2 + \hat{\text{df}}_{A_T} \frac{\log(n)}{n}\right) \rightarrow 1, \tag{27}$$

$$P\left(\inf_{(\lambda, \nu) \in \Omega_-} \text{BIC}(\lambda, \nu) > \text{BIC}(\lambda_n^*, \nu_n^*)\right) \rightarrow 1, \tag{28}$$

$$P\left(\inf_{(\lambda, \nu) \in \Omega_+} \text{BIC}(\lambda, \nu) > \text{BIC}(\lambda_n^*, \nu_n^*)\right) \rightarrow 1, \tag{29}$$

where  $\hat{\text{df}}_{A_T}$  is the degrees of freedom of  $\hat{\mathbf{Y}}_{A_T} = \mathbf{X}\hat{\mathbf{b}}_{A_T} + \mathbf{Z}\hat{\mathbf{u}}_{A_T}$ .

For (27), it follows from theorem 1 that  $P(\hat{\beta}_{\lambda_n^*, \nu_n^*} = \hat{b}_{A_T}) \rightarrow 1$ , as  $n \rightarrow \infty$ . By (27), (B.1) and (B.2) with  $(\lambda, \nu) \in \Omega_-$ , we have as  $n \rightarrow \infty$ ,

$$\begin{aligned} & \text{BIC}(\lambda_n^*, \nu_n^*) \\ &= \log(\hat{\sigma}_{\hat{A}_{\lambda_n^*, \nu_n^*}}^2) + \frac{(\mathbf{Y} - \mathbf{X}\hat{\beta}_{\lambda_n^*, \nu_n^*})' \mathbf{V}^{*-2} (\mathbf{Y} - \mathbf{X}\hat{\beta}_{\lambda_n^*, \nu_n^*})}{(\mathbf{Y} - \mathbf{X}\hat{\beta}_{\lambda_n^*, \nu_n^*})' \mathbf{V}^{*-1} (\mathbf{Y} - \mathbf{X}\hat{\beta}_{\lambda_n^*, \nu_n^*})} + o_p(1) \\ &\geq \log(\hat{\sigma}_{\hat{A}_{\lambda_n^*, \nu_n^*}}^2) + \frac{(\mathbf{Y} - \mathbf{X}\hat{\beta}_{\lambda_n^*, \nu_n^*})' \mathbf{V}^{*-1/2} (\mathbf{I} - \mathbf{Z}(\mathbf{Z}'\mathbf{Z} + \mathbf{G}^{-1})^{-1} \mathbf{Z}') \mathbf{V}^{*-1/2} (\mathbf{Y} - \mathbf{X}\hat{\beta}_{\lambda_n^*, \nu_n^*})}{(\mathbf{Y} - \mathbf{X}\hat{\beta}_{\lambda_n^*, \nu_n^*})' \mathbf{V}^{*-1} (\mathbf{Y} - \mathbf{X}\hat{\beta}_{\lambda_n^*, \nu_n^*})} \\ &\quad + o_p(1) \\ &\rightarrow \log(\sigma_{A_T}^2) + 1 + o_p(1) \end{aligned}$$

and

$$\begin{aligned}
& \text{BIC}(\lambda, \nu) \\
&= \log(\hat{\sigma}_{\lambda, \nu}^2) + \frac{1}{n\hat{\sigma}_{\lambda, \nu}^2} (\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}_{\lambda, \nu} - \mathbf{Z}\hat{\mathbf{u}}_{\lambda, \nu})' (\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}_{\lambda, \nu} - \mathbf{Z}\hat{\mathbf{u}}_{\lambda, \nu}) + o_p(1) \\
&\geq \log(\hat{\sigma}_{\hat{A}_{\lambda, \nu}}^2) + \frac{(\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}_{\lambda, \nu})' \mathbf{V}^{*-2} (\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}_{\lambda, \nu})}{(\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}_{\lambda, \nu})' \mathbf{V}^{*-1} (\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}_{\lambda, \nu})} + o_p(1) \\
&\geq \log(\hat{\sigma}_{\hat{A}_{\lambda, \nu}}^2) + \frac{(\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}_{\lambda, \nu})' \mathbf{V}^{*-1/2} (\mathbf{I} - \mathbf{Z}(\mathbf{Z}'\mathbf{Z} + \mathbf{G}^{-1})^{-1} \mathbf{Z}') \mathbf{V}^{*-1/2} (\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}_{\lambda, \nu})}{(\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}_{\lambda, \nu})' \mathbf{V}^{*-1} (\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}_{\lambda, \nu})} \\
&\quad + o_p(1) \\
&\rightarrow \min_{A: A \cap A_T \neq A_T} \log(\sigma_A^2) + 1 > \log(\sigma_{A_T}^2) + 1,
\end{aligned}$$

which imply (28). It remains to show (29). Under (B.3),

$$\left\{ \hat{\text{df}}_{\hat{\mathbf{Y}}}(\lambda, \nu) - \hat{\text{df}}_{\hat{\mathbf{Y}}}(\lambda_n^*, \nu_n^*) \right\} - \left\{ \hat{\text{df}}_{\hat{A}_{\lambda, \nu}} - \hat{\text{df}}_{A_T} \right\} \rightarrow 0 \text{ in probability,}$$

Therefore, as  $n \rightarrow \infty$ ,

$$\begin{aligned}
& n \left( \text{BIC}(\lambda, \nu) - \text{BIC}(\lambda_n^*, \nu_n^*) \right) \\
&= n \log \left( \frac{\hat{\sigma}_{\lambda, \nu}^2}{\hat{\sigma}_{\lambda_n^*, \nu_n^*}^2} \right) + \frac{\hat{\sigma}_{A_T}^2 \|\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}_{\lambda, \nu} - \mathbf{Z}\hat{\mathbf{u}}_{\lambda, \nu}\|^2}{\hat{\sigma}_{\lambda, \nu}^2 \|\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}_{A_T} - \mathbf{Z}\hat{\mathbf{u}}_{A_T}\|^2} \\
&\quad + \left\{ \hat{\text{df}}_{\hat{\mathbf{Y}}}(\mathbf{G}, \lambda, \nu) - \hat{\text{df}}_{\hat{\mathbf{Y}}}(\mathbf{G}, \lambda_n^*, \nu_n^*) \right\} \log(n) \\
&\rightarrow \hat{\sigma}_{A_T}^{-2} n (\hat{\sigma}_{\lambda, \nu}^2 - \hat{\sigma}_{A_T}^2) + o_p(1) + 1 + \left\{ \hat{\text{df}}_{\hat{A}_{\lambda, \nu}} - \hat{\text{df}}_{A_T} \right\} \log(n).
\end{aligned}$$

It follows that

$$\begin{aligned}
\inf_{(\lambda, \nu) \in \Omega_+} n \left( \text{BIC}(\lambda, \nu) - \text{BIC}(\lambda_n^*, \nu_n^*) \right) &\geq \hat{\sigma}_{A_T}^{-2} \min_{A: A \supseteq A_T} n (\hat{\sigma}_A^2 - \hat{\sigma}_{A_T}^2) + 1 \\
&\quad + (\hat{\text{df}}_{\hat{A}_{\lambda, \nu}} - \hat{\text{df}}_{A_T}) \log(n) + o_p(1). \quad (30)
\end{aligned}$$

Note that  $n(\hat{\sigma}_{A_T}^2 - \hat{\sigma}_A^2) = O_p(1)$  follows non-central chi-square distribution. By (17) and (18),

$$\hat{\text{df}}_A = \text{tr} \left[ \left\{ \mathbf{I} - \mathbf{Z}(\mathbf{Z}'\mathbf{Z} + \mathbf{G}^{-1})^{-1} \mathbf{Z}' \right\} \mathbf{V}^{-1} \mathbf{H}_A + \mathbf{Z}(\mathbf{Z}'\mathbf{Z} + \mathbf{G}^{-1})^{-1} \mathbf{Z}' \right]$$

and

$$\hat{\text{df}}_{\hat{A}_{\lambda, \nu}} - \hat{\text{df}}_{A_T} = \text{tr} \left[ \left\{ \mathbf{I} - \mathbf{Z}(\mathbf{Z}'\mathbf{Z} + \mathbf{G}^{-1})^{-1} \mathbf{Z}' \right\} \mathbf{V}^{-1} \left\{ \mathbf{H}_{\hat{A}_{\lambda, \nu}} - \mathbf{H}_{A_T} \right\} \right],$$

where  $\mathbf{H}_A = \mathbf{X}_A (\mathbf{X}'_A \mathbf{V}^{-1} \mathbf{X}_A)^{-1} \mathbf{X}'_A$ .

It can be shown that  $\mathbf{I} - \mathbf{Z}(\mathbf{Z}'\mathbf{Z} + \mathbf{G}^{-1})^{-1} \mathbf{Z}'$  and  $\mathbf{H}_{\hat{A}_{\lambda, \nu}} - \mathbf{H}_{A_T}$  are strictly positive definite, therefore  $\hat{\text{df}}_{\hat{A}_{\lambda, \nu}} - \hat{\text{df}}_{A_T} > 0$  for  $(\lambda, \nu) \in \Omega_+$ . The righthand side of (30) diverges to  $+\infty$  as  $n \rightarrow \infty$ , and hence (29) is satisfied. This completes the proof.  $\square$

## References

- BONDELL, H. D., KRISHNA, A. AND GHOSH, S. K. (2010) Joint variable selection for fixed and random effects in linear mixed-effects models. *Biometrics*, **66**, 1069–1077. [MR2758494](#)
- DE BOOR, C. (2001) *A practical guide to splines*. New York: Springer Verlag. [MR1900298](#)
- CHEN, Z. AND DUNSON, D. B. (2003) Random effects selection in linear mixed models. *Biometrics*, **59**, 762–769. [MR2025100](#)
- DIGGLE, P., HEAGERTY, P., LIANG, K. AND ZEGER, S. (2002) *Analysis of longitudinal data*. USA: Oxford University Press. [MR2049007](#)
- FAHRMEIR, L. AND LANG, S. (2001) Bayesian inference for generalized additive mixed models based on markov random field priors. *Applied Statistics*, 201–220. [MR1833273](#)
- HARVILLE, D. A. (1974) Bayesian inference for variance components using only error contrasts. *Biometrika*, **61**, 383–385. [MR0368279](#)
- HEDEKER, D. AND GIBBONS, R. (2006) *Longitudinal data analysis*. New York: Wiley. [MR2284230](#)
- HUANG, J., HOROWITZ, J. AND WEI, F. (2010) Variable selection in nonparametric additive models. *The Annals of Statistics*, **38**, 2282–2313. [MR2676890](#)
- KINNEY, S. K. AND DUNSON, D. B. (2007) Fixed and random effects selection in linear and logistic models. *Biometrics*, **63**, 690–698. [MR2395705](#)
- LAIRD, N. AND WARE, J. (1982) Random-effects models for longitudinal data. *Biometrics*, **38**, 963–974.
- LIANG, K. AND ZEGER, S. (1986) Longitudinal data analysis using generalized linear models. *Biometrika*, **73**, 13. [MR0836430](#)
- LIN, X. AND ZHANG, D. (1999) Inference in generalized additive mixed models by using smoothing splines. *Journal of the Royal Statistical Society Series B*, **61**, 381–400. [MR1680318](#)
- LIN, Y. AND ZHANG, H. (2006) Component selection and smoothing in multivariate nonparametric regression. *The Annals of Statistics*, **34**, 2272. [MR2291500](#)
- LINDSTROM, M. AND BATES, D. (1988) Newton-Raphson and EM algorithms for linear mixed-effects models for repeated-measures data. *Journal of the American Statistical Association*, 1014–1022. [MR0997577](#)
- MEIER, L., VAN DE GEER, S. AND BUHLMANN, P. (2008) The group lasso for logistic regression. *Journal of the Royal Statistical Society Series B*, **70**, 53. [MR2412631](#)
- MEIER, L., VAN DE GEER, S., BUHLMANN, P. AND ZURICH, E. (2009) High-dimensional additive modeling. *The Annals of Statistics*, **37**, 3779–3821. [MR2572443](#)
- PU, W. AND NIU, X.-F. (2006) Selecting mixed-effects models based on a generalized information criterion. *JMA*, **97**, 733–758. [MR2236499](#)
- RAVIKUMAR, P., LAFFERTY, J., LIU, H. AND WASSERMAN, L. (2009) Sparse additive models. *Journal of the Royal Statistical Society Series B*, **71**, 1009–1030. [MR2750255](#)



- REISBY, N., GRAM, L., BECH, P., NAGY, A., PETERSEN, G., ORTMANN, J., IBSEN, I., DENCKER, S., JACOBSEN, O., KRAUTWALD, O. *et al.* (1977) Imipramine: clinical effects and pharmacokinetic variability. *Psychopharmacology*, **54**, 263–272.
- RUPPERT, D., WAND, M. AND CARROLL, R. (2003) *Semiparametric regression*. New York: Cambridge University Press. [MR1998720](#)
- SHEN, X. AND YE, J. (2002) Adaptive model selection. *Journal of the American Statistical Association*, **97**, 210–221. [MR1947281](#)
- STONE, C. (1985) Additive regression and other nonparametric models. *The Annals of Statistics*, **13**, 689–705. [MR0790566](#)
- STONE, C. (1986) The dimensionality reduction principle for generalized additive models. *The Annals of Statistics*, 590–606. [MR0840516](#)
- TIBSHIRANI, R. (1996) Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society Series B*, **58**, 267–288. [MR1379242](#)
- VAN DER VAART, A. AND WELLNER, J. (1996) *Weak convergence and empirical processes: with applications to statistics*. New York: Springer Verlag. [MR1385671](#)
- VAN DE GEER, S. (2008) High-dimensional generalized linear models and the lasso. *The Annals of Statistics*, **36**, 614. [MR2396809](#)
- WAND, M. (2003) Smoothing and mixed models. *Computational Statistics*, **18**, 223–250.
- WANG, H., LI, R. AND TSAI, C. (2007a) Tuning parameter selectors for the smoothly clipped absolute deviation method. *Biometrika*, **94**, 553–568. [MR2410008](#)
- WANG, L., CHEN, G. AND LI, H. (2007b) Group scad regression analysis for microarray time course gene expression data. *Bioinformatics*, **23**, 1486.
- WEI, F. AND HUANG, J. (2008) Consistent group selection in high-dimensional linear regression. *Tech. rep.*, Department of Statistics and Actuarial Science, University of Iowa.
- YUAN, M. AND LIN, Y. (2006) Model selection and estimation in regression with grouped variables. *Journal of the Royal Statistical Society Series B*, **68**, 49–67. [MR2212574](#)
- ZEGER, S. AND LIANG, K. (1986) Longitudinal data analysis for discrete and continuous outcomes. *Biometrics*, **42**, 121–130.
- ZHANG, D., LIN, X., RAZ, J. AND SOWERS, M. (1998) Semiparametric stochastic mixed models for longitudinal data. *Journal of the American Statistical Association*, **93**. [MR1631369](#)
- ZHOU, S., SHEN, X. AND WOLFE, D. (1998) Local asymptotics for regression splines and confidence regions. *The Annals of Statistics*, **26**, 1760–1782. [MR1673277](#)
- ZOU, H. (2006) The adaptive lasso and its oracle properties. *Journal of the American Statistical Association*, **101**, 1418–1429. [MR2279469](#)