# ESTIMATION OF SPARSE FUNCTIONAL ADDITIVE MODELS WITH ADAPTIVE GROUP LASSO

Peijun Sang[1], Liangliang Wang[2] and Jiguo Cao[2]

[1]*University of Waterloo* and [2]*Simon Fraser University*

*Abstract:* We study a flexible model to address the lack of fit in conventional functional linear regression models. This model, called the sparse functional additive model, is used to characterize the relationship between a functional predictor and a scalar response of interest. The effect of the functional predictor is represented in a nonparametric additive form, where the arguments are the scaled functional principal component scores. Component selection and smoothing are considered when fitting the model in order to reduce the variability and enhance the prediction accuracy, while providing an adequate fit. To achieve these goals, we propose using the adaptive group LASSO method to select relevant components and smoothing splines and, thus, obtain a smoother estimate of those relevant components. Simulation studies show that the proposed estimation method compares favorably with conventional methods in terms of prediction accuracy and component selection. Furthermore, the advantages of our estimation method are demonstrated using two real-data examples.

*Key words and phrases:* Functional data analysis, functional linear model, functional principal component analysis, group LASSO, smoothing spline.

## 1. Introduction

Functional data analysis has become an important tool for dealing with data collected over multiple time points, spatial locations, or other continua. A fundamental problem in functional data analysis is how to model the relationship between a scalar response of interest and a functional predictor. For instance, the Tecator data (see Section 5.1) measure 240 meat samples, each of which has a spectrum of absorbance and contains water, fat, and protein. Researchers have investigated how to use the spectrum of absorbance, which can be treated as a functional predictor, to predict one of the three contents. A functional linear regression (FLR) is a conventional and interpretable model for predicting a scalar response from a functional predictor. It has many interesting applications. For instance, Ainsworth, Routledge and Cao (2011) applied an FLR to explore the effect of river flow on the decline of sockeye salmon. Luo et al. (2013) applied an

FLR to investigate the time-varying intensity of ward admission and its effect on to emergency department access block.

In an FLR, the relationship between a scalar response and a functional predictor is modeled in a linear form. Hence, the key to fitting an FLR is to estimate the coefficient function of the functional predictor. There has been extensive research to address this problem. For example, Müller and Stadtmüller (2005) represented the coefficient function in terms of Fourier basis functions or the eigenfunctions of the estimated covariance function of the functional predictor. Then, the coefficients of the Fourier basis functions are obtained by solving a functional estimating equation. Ramsay and Silverman (2005) suggested using spline basis functions to represent the coefficient function. Then, they solve a regularized regression problem, in which the roughness of the spline representation is penalized to obtain a smooth estimate of the coefficient function. Lin et al. (2017) proposed a local sparse estimator for the coefficient function to enhance the interpretability of FLRs. Liu, Wang and Cao (2017) added a random effect on the coefficient function when repeated measurements are available on multiple subjects. A comprehensive introduction to FLRs can be found in Horváth and Kokoszka (2012) and Morris (2015).

Although the aforementioned studies have proposed various estimation methods that can be used to fit an FLR model, and have established some appealing properties of the corresponding estimators, in practice, applications of FLRs can be restricted, owing to its simple linear form. Similarly to the multiple linear model, which in some cases may not adequately describe the relationship between a scalar response and scalar covariates, an FLR may suffer from inadequate flexibility in terms of modelling the relationship between a scalar response and a functional predictor. This phenomenon has been noted by many researchers. For instance, Yao and Müller (2010) extended the FLR model to the case when the scalar response depends on a polynomial of the functional predictor, focusing mainly on the quadratic case. Chen et al. (2011) used a nonparametric link to connect the scalar response and the functional linear form. A class of flexible functional nonlinear regression models has been proposed by Müller, Wu and Yao (2013), who use continuously additive models to characterize the relationship between a functional predictor and a scalar response. Nonlinear and/or nonparametric functional regression models can somewhat address the issue of an inadequate fit caused by an FLR (see Chen et al. (2011), Müller, Wu and Yao (2013), Müller and Yao (2008)). However, these models have other disadvantages such as over-flexibility and a lack of stability (Zhu, Yao and Zhang (2014)).

Reiss et al. (2017) summarized some of main approaches used to regress a scalar response on a functional predictor. We propose a functional regression model that achieves a satisfactory tradeoff between flexibility and simplicity.

Zhu, Yao and Zhang (2014) proposed an extended functional additive model, in which the scalar response of interest depends on a transformation of the leading functional principal component (FPC) scores. They assumed that some additive components were vanishing, and that the nonvanishing components were smooth functions, for the sake of simplicity and interpretability, while retaining flexibility. To achieve this goal, they adopted the regularization scheme of the component selection and smoothing operator (COSSO) proposed by Lin and Zhang (2006), which can select and smooth components simultaneously. This model achieves a better tradeoff between flexibility and simplicity than many other functional regression models do. However, the estimation procedure seems to suffer from several drawbacks. First, only estimation consistency is guaranteed for the proposed estimator. Whether selection consistency holds for this estimator remains an open question. Another drawback is associated with computational complexity. As noted by Zhang and Lin (2006), when a full basis is employed, the complexity of the algorithm is $O(n^3)$, where $n$ is the sample size. To reduce the computational burden, Zhang and Lin (2006) suggested using a subset basis algorithm instead, which was computationally much more efficient than the full basis algorithm. Zhu, Yao and Zhang (2014) seemed to ignore this computational issue when implementing COSSO to fit the proposed model. The computational complexity is demonstrated in simulation studies.

To overcome the drawbacks of the method proposed by Zhu, Yao and Zhang (2014), we propose a method for estimating extended functional additive models. In contrast to representing nonparametric additive components in the framework of RKHS (Zhu, Yao and Zhang (2014)), we use B-spline basis functions to represent these components, which are easier to understand and implement. Then, selecting nonzero components is equivalent to selecting the nonzero coefficients of the B-spline basis functions. The group LASSO method (Yuan and Lin (2006)) has been shown to perform well when selecting grouped variables for accurate prediction, in both theory and application. Because an additive component corresponds to a vector of coefficients, which can be treated as a group of variables, we employ the group LASSO method to select nonzero vectors of coefficients. The adaptive group LASSO method is then applied to allow for variation in the shrinkages of the vectors of the coefficients. This modification yields a more accurate estimate of the coefficient vectors, which then leads to a better estimate

for the additive components. This method enables us to achieve our goal of obtaining a parsimonious model via component selection.

Nevertheless, the estimated nonzero components can be wiggly, because we represent the additive components using a large number of B-spline basis functions. This may impair the predictive performance, as shown in the simulation studies in Section 4. Thus, we suggest refining the selected components using smoothing splines. This extra smoothing step improves the prediction accuracy of the estimator obtained from the adaptive group LASSO, as shown in our simulation studies.

This study makes three main contributions to the literature. First, compared with traditional FLR models, our proposed model provides a better tradeoff between flexibility and simplicity when modeling the effect of a functional predictor. By selecting and smoothing nonzero components, our proposed method obtains an estimator that has better prediction accuracy. Second, unlike the COSSO regularization scheme adopted in Zhu, Yao and Zhang (2014), we employ a group LASSO to select components, and use the smoothing spline method to smooth nonzero components. As a result, our proposed estimation method is easy to understand and implement. Last, but not least, we provide both theoretical and empirical examples of the selection consistency and estimation consistency of our proposed estimator; in contrast, Zhu, Yao and Zhang (2014) provided only a theoretical proof of the estimation consistency of their estimator.

The remainder of this paper is organized as follows. Section 2 introduces a sparse functional additive model, and our method for estimating the additive components in the model. Section 3 establishes the selection consistency and the estimation consistency of our proposed estimator. The finite-sample performance of the estimator is investigated empirically in Section 4, where we conduct simulation studies to compare our proposed estimator with other conventional methods. In Section 5, our method is demonstrated by analyzing two real-data examples. Section 6 concludes the paper. The procedures used to estimate the FPC scores, proofs of the main results in Section 3, and additional empirical studies are provided in the online Supplementary Material.

## 2. Model and Estimation Method

### 2.1. Sparse functional additive model

Suppose that $\{X_i(t), y_i\}_{i=1}^n$ are independent and identically distributed (i.i.d.) observations from $\{X(t), Y\}$, where $X(t)$ is a random function, and $Y$ is a scalar

random variable. We assume $X(t)$ is a square integrable stochastic process over a compact interval $\mathcal{I} = [0, T]$; that is, $\mathrm{E}\left\{\int_{\mathcal{I}} X^2(t)dt\right\} < \infty$. Let $m(t)$ and $G(s,t)$ denote the mean function and covariance function of $X(t)$, respectively. According to Mercer's theorem, $G(s,t)$ can be represented as $G(s,t) = \sum_{k=1}^{\infty} \lambda_k \phi_k(s) \phi_k(t)$, where $\lambda_k$ is a nonnegative eigenvalue, and $\phi_k(t)$ is the corresponding eigenfunction. For the sake of identifiability, we postulate that $\lambda_1 \geq \lambda_2 \geq \cdots \geq 0$. Additionally, $\{\phi_k\}_{k=1}^{\infty}$ is assumed to be a complete orthonormal basis of the space $L^2(\mathcal{I})$, the collection of all square integrable functions on $\mathcal{I}$. Then, the stochastic process $X(t)$ admits the Karhunen-Loève expansion:

$$X(t) = m(t) + \sum_{k=1}^{\infty} \xi_k \phi_k(t), \tag{2.1}$$

where $\xi_k = \int_{\mathcal{I}} (X(t) - m(t)) \phi_k(t) dt$, for $k = 1, \ldots$, is called the $k$th FPC score. The FPC score satisfies $\mathrm{E}(\xi_k \xi_{k'}) = \lambda_k$ if $k = k'$, and is zero otherwise.

In an FLR, $Y$ is treated as the response and $X(t)$ is the functional predictor. Furthermore, the relationship between $Y$ and $X(t)$ is modeled in a linear form:

$$y_i = \int_{\mathcal{I}} X_i(t) b(t) dt + \epsilon_i,$$

where $\epsilon_i$ denotes a random error with mean zero and variance $\sigma_\epsilon^2$. Given the representation of $X(t)$ in (2.1), we have $y_i = a + \sum_{k=1}^{\infty} b_k \xi_{ik} + \epsilon_i$, where $a = \int_{\mathcal{I}} m(t) b(t) dt$, $\xi_{ik}$ denotes the $k$th FPC score of $X_i(t)$, and $b_k = \int_{\mathcal{I}} \phi_k(t) b(t) dt$, for $k \geq 1$. To address the curse of dimensionality, a truncated model is usually adopted, such that $Y$ depends only on the first $d$ FPC scores. In other words, we get a truncated linear model: $y_i = a + \sum_{j=1}^{d} b_j \xi_{ij} + \epsilon_i$. In practice, $d$ is chosen as the smallest number of FPCs that can explain over 99.9% of the total variability of the functional predictor $X(t)$. As noted by Zhu, Yao and Zhang (2014), this choice can, to some extent, circumvent neglecting those FPC scores that play a negligible role in capturing the variability of the functional predictor, but that are relevant to predicting the response. This truncated model is slightly restrictive, because an explicit parametric form is assumed between the response and the leading FPC scores. The linearity assumption is likely to be violated in most practical scenarios.

Based on the work of Hastie and Tibshirani (1986), and the fact that the $\xi_j$s are mutually uncorrelated, a nonparametric functional additive model was proposed by Müller and Yao (2008) to describe the relationship between the

response and the first $d$ FPC scores,

$$y_i = a + \sum_{j=1}^{d} f_j(\xi_{ij}) + \epsilon_i, \qquad (2.2)$$

where we call $f_j$ the $j$th component in the nonparametric functional additive model.

FPC scores usually cannot be observed directly. Therefore, we first need to estimate the FPC scores from the observed functional data, which may be subject to measurement errors. We assume that $W_{ij} = X_i(t_{ij}) + e_{ij}$, where $W_{ij}$ denotes the observation of the process $X_i(t)$, made at time point $t_{ij}$, for $j = 1, \ldots, N_i$, $i = 1, \ldots, n$. Furthermore, $e_{ij}$ denotes a measurement error, and is assumed to be independent of $X_i(t)$. The functional principal component analysis (FPCA) is implemented to estimate the FPC scores, denoted by $\hat{\xi}_{ij}$. The details of this procedure can be found in the Supplementary Material.

We first scale the FPC scores to $[0, 1]$ using a transformation function $F$. One possible strategy is to apply the cumulative distribution function (cdf) $F(z|\lambda_j)$ of the Normal$(0, \lambda_j)$ on $\xi_j$, where $\lambda_j$ is the eigenvalue of the covariance function $G(s, t)$, and $\lambda_j = Var(\xi_j)$. We define $\zeta_j$ as the $j$th scaled FPC score: $\zeta_j = F(\xi_j|\lambda_j)$, for $j = 1, \ldots, d$. Then, the estimated scaled FPC scores are given as $\hat{\zeta}_{ij} = F(\hat{\xi}_{ij}|\hat{\lambda}_j)$, for $j = 1, \ldots, d$, $i = 1, \ldots, n$. Assumption B in Section 3 provides conditions on the transformation function $F$ that are more general. We still use $f_j$ for the $j$th component in the nonparametric functional additive model when $\xi_j$ is replaced by $\zeta_j$.

The nonparametric functional additive model (2.2) can now be expressed as

$$y_i = a + \sum_{j=1}^{d} f_j(\zeta_{ij}) + \epsilon_i. \qquad (2.3)$$

To make the model identifiable, we assume that $\mathrm{E}\{f_j(\zeta_j)\} = 0$, for $j = 1, \ldots, d$. Models with a parsimonious structure are preferable, in practice. Thus, we assume that some components, $f_j$ are vanishing, and that the remainder of the components are nonzero and smooth. Model (2.3) is called a sparse functional additive model in this article.

B-spline functions, owing to their nice properties (De Boor (2001)), are widely used to estimate unknown functions (see Stone (1985, 1986); Huang, Horowitz and Wei (2010), etc). In this study, we employ B-spline functions to es-

timate the additive components in Model (2.3). We begin with a brief overview of B-splines. For more information, see De Boor (2001). Let $0 = \tau_0 < \tau_1 < \cdots < \tau_{L_n} < \tau_{L_n+1} = 1$ be the breakpoints that separate the interval $[0, 1]$ into $L_n + 1$ subintervals. We assume that $L_n = O(n^\alpha)$, where $0 < \alpha < 0.5$, and define $\delta_n = \max_{0 \le m \le L_n} |\tau_{m+1} - \tau_m| = O(n^{-\alpha})$. Let $c_1$ be a constant, independent of $n$, such that $\delta_n < c_1 \min_{0 \le m \le L_n} |\tau_{m+1} - \tau_m|$. Let $\mathscr{S}_n$ be the space of polynomial splines of order $l$, which is one more than the degree of polynomials, on $[0, 1]$ consisting of functions $s$ satisfying the following: (i) $s$ is a polynomial of order $l$ at each subinterval $[\tau_m, \tau_{m+1}]$, for $m = 0, \ldots, L_n$; and (ii) for $0 \le l^\star \le l - 2$, the $l^\star$th-order derivative of $s$ is continuous on $[0, 1]$. Then, there exist $m_n = L_n + l$ normalized B-spline basis functions $\{B_k, 1 \le k \le m_n\}$, bounded by zero and one on $[0, 1]$, such that any $f \in \mathscr{S}_n$ can be written as

$$f_j(x) = \sum_{k=1}^{m_n} \beta_{jk} B_k(x), \quad j = 1, 2, \ldots, d, \tag{2.4}$$

where $\boldsymbol{\beta}_j = (\beta_{j1}, \ldots, \beta_{jm_n})'$ is the spline coefficient vector. Now, selecting nonzero components $f_j(\cdot)$ for Model (2.3) amounts to selecting nonzero $\boldsymbol{\beta}_j$.

## 2.2. Group LASSO

Accounting for the fact that $\mathrm{E}\{f_j(\zeta_j)\} = 0$, for $j = 1, \ldots, d$, we define $\psi_{jk}(x) = B_k(x) - (1/n)\sum_{i=1}^n B_k(\hat{\zeta}_{ij})$, for $k = 1, \ldots, m_n$, $j = 1, \ldots, d$. For brevity, $\psi_{jk}(x)$ is denoted by $\psi_k(x)$, without causing any confusion. Thus, $\sum_{i=1}^n \psi_k(\hat{\zeta}_{ij}) = 0$, for $j = 1, \ldots, d$. The estimated intercept in Model (2.3) is given as $\bar{y} = \frac{1}{n}\sum_{i=1}^n y_i$. Let $\boldsymbol{Z}_{ij} = (\psi_1(\hat{\zeta}_{ij}), \ldots, \psi_{m_n}(\hat{\zeta}_{ij}))^T$, $\boldsymbol{Z}_j = (\boldsymbol{Z}_{1j}, \ldots, \boldsymbol{Z}_{nj})^T$, and $\boldsymbol{Z} = (\boldsymbol{Z}_1, \ldots, \boldsymbol{Z}_d)$. Similarly, define $\boldsymbol{\beta} = (\boldsymbol{\beta}_1^T, \ldots, \boldsymbol{\beta}_d^T)^T$, where $\boldsymbol{\beta}_j = (\beta_{j1}, \ldots, \beta_{jm_n})^T$, and $\boldsymbol{y} = (y_1 - \bar{y}, \ldots, y_n - \bar{y})^T$. Nonzero $\boldsymbol{\beta}_j$ in Model (2.3) can be selected and estimated using the group LASSO (Yuan and Lin (2006)), in which the corresponding estimate $\tilde{\boldsymbol{\beta}}$ minimizes

$$D_1(\boldsymbol{\beta}) = (\boldsymbol{y} - \boldsymbol{Z}\boldsymbol{\beta})^T(\boldsymbol{y} - \boldsymbol{Z}\boldsymbol{\beta}) + \lambda_1 \sum_{j=1}^d ||\boldsymbol{\beta}_j||_2. \tag{2.5}$$

In (2.5), the positive tuning parameter $\lambda_1$ determines the magnitude of the shrinkage, and $|| \cdot ||_2$ denotes the Euclidean norm of a vector in $\mathbb{R}^{m_n}$. If $\tilde{\boldsymbol{\beta}}_j = (\tilde{\beta}_{j1}, \ldots, \tilde{\beta}_{jm_n})^T$, then the corresponding estimate of $f_j$ is denoted by $\tilde{f}_j$, which is equal to $\sum_{k=1}^{m_n} \tilde{\beta}_{jk} \psi_k(x)$. Cross-validation is employed to choose an "optimal" $\lambda_1$, which is chosen to minimize the cross-validation error.

## 2.3. Adaptive group LASSO

The group LASSO method penalizes each $\boldsymbol{\beta}_j$ equally in (2.5), which may not be an optimal treatment. Thus, to account for different impacts on $\zeta_j$, we propose an adaptive group LASSO method, which is similar in spirit to the adaptive LASSO method proposed by Zou (2006). More explicitly, we introduce a weight vector $(w_1, \ldots, w_d)$, which allows each $\boldsymbol{\beta}_j$ to have its own shrinkage value. Given $\tilde{\boldsymbol{\beta}}$, estimated using the group LASSO, for $j = 1, \ldots, d$, $w_j$ is set as $||\tilde{\boldsymbol{\beta}}_j||_2^{-1}$ if $||\tilde{\boldsymbol{\beta}}_j||_2 > 0$, and $\infty$ otherwise. Then, the adaptive group LASSO estimate of $\boldsymbol{\beta}$, denoted by $\widehat{\boldsymbol{\beta}}$, is obtained by minimizing

$$D_2(\boldsymbol{\beta}) = (\boldsymbol{y} - \boldsymbol{Z}\boldsymbol{\beta})^T(\boldsymbol{y} - \boldsymbol{Z}\boldsymbol{\beta}) + \lambda_2 \sum_{j=1}^d w_j ||\boldsymbol{\beta}_j||_2, \qquad (2.6)$$

where $\lambda_2$ denotes a penalty parameter that can be determined by cross-validation. Then, the corresponding estimate of $f_j(x)$, denoted by $\hat{f}_j(x)$, can be represented in terms of $\boldsymbol{\psi}_j(x) = (\psi_{j1}(x), \ldots, \psi_{jm_n}(x))^T$; that is, $\hat{f}_j(x) = \widehat{\boldsymbol{\beta}}_j^T \boldsymbol{\psi}_j(x)$, for $j = 1, \ldots, d$. If $\widehat{\boldsymbol{\beta}}_j = \boldsymbol{0}$ for some $j$, then the estimate $\hat{f}_j$ is also zero.

## 2.4. Smoothing spline method

When a large number of B-spline basis functions are employed to estimate $f_j$, then the adaptive group LASSO estimate may be wiggly, in which case, further smoothing of the nonzero estimates obtained from the adaptive group LASSO is required. This concern is also discussed in Wu et al. (2014). To allow for different roughness penalties for the nonzero components, we propose a smoothing spline method. The weight is defined as $w_j = ||\widehat{\boldsymbol{\beta}}_j||_2^{-1}$, where $j \in S$, and $S = \{j : \widehat{\boldsymbol{\beta}}_j \neq \boldsymbol{0}\}$ is the set of nonzero components. In particular, the updated estimate of $\boldsymbol{\beta}_j$ is obtained from the smoothing spline method by minimizing

$$D_3(\boldsymbol{\beta}) = \left(\boldsymbol{y} - \sum_{j \in S} \boldsymbol{Z}_j \boldsymbol{\beta}_j\right)^T \left(\boldsymbol{y} - \sum_{j \in S} \boldsymbol{Z}_j \boldsymbol{\beta}_j\right) + \lambda_3 \sum_{j \in S} w_j \int_0^1 \{f_j''(\zeta_j)\}^2 d\zeta_j, \quad (2.7)$$

where $\lambda_3$ denotes the smoothing parameter. The roughness penalty term $\int_0^1 \{f_j''(\zeta_j)\}^2 d\zeta_j = \boldsymbol{\beta}_j^T \boldsymbol{Q}_j \boldsymbol{\beta}_j$, where $\boldsymbol{Q}_j$ is an $m_n \times m_n$ penalty matrix, with $pq$th element $\boldsymbol{Q}_j^{pq} = \int_0^1 B_p''(\zeta_j) B_q''(\zeta_j) d\zeta_j$. When the second derivative of $f_j(\zeta_j)$ does not exist, the penalty matrix $\boldsymbol{Q}_j$ can be replaced by the difference matrix introduced by Eilers and Marx (1996). Minimizing (2.7) is equivalent to

a classical smoothing spline problem, except that there is a weight vector in this problem. Let $\sum_{j \in S} \boldsymbol{Z}_j \boldsymbol{\beta}_j = \boldsymbol{Z}_S \boldsymbol{\beta}$, where $\boldsymbol{Z}_S = (\boldsymbol{Z}_{i_1}, \ldots, \boldsymbol{Z}_{i_{|S|}}) \in \mathbb{R}^{n \times m_n |S|}$, $i_1, \ldots, i_{|S|}$ are all elements of $S$, and $|S|$ denotes the cardinality of the set $S$. Let $\boldsymbol{Q} = \mathrm{diag}(w_{i_1} \boldsymbol{Q}_{i_1}, \ldots, w_{i_{|S|}} \boldsymbol{Q}_{i_{|S|}})$. Then, the estimate of $\boldsymbol{\beta}$, still denoted by $\widehat{\boldsymbol{\beta}}$, is given as $\widehat{\boldsymbol{\beta}} = (\boldsymbol{Z}_S^T \boldsymbol{Z}_S + \lambda_3 \boldsymbol{Q})^{-1} \boldsymbol{Z}_S^T \boldsymbol{y}$. The corresponding estimate of $f_j$ is $\hat{f}_j = \widehat{\boldsymbol{\beta}}_j^T \boldsymbol{\psi}_j(x)$, for $j \in S$.

The smoothing parameter $\lambda_3$ can be determined by the generalized cross-validation (GCV) measure. For a given $\lambda_3$, the corresponding measure can be expressed as

$$\mathrm{GCV}(\lambda_3) = \frac{n \cdot \mathrm{SSE}}{(n - df(\lambda_3))^2},$$

where $\mathrm{SSE} = (\boldsymbol{y} - \boldsymbol{Z}_S \widehat{\boldsymbol{\beta}})^T (\boldsymbol{y} - \boldsymbol{Z}_S \widehat{\boldsymbol{\beta}})$ and $df(\lambda_3) = \mathrm{trace}(\boldsymbol{Z}_S (\boldsymbol{Z}_S^T \boldsymbol{Z}_S + \lambda_3 \boldsymbol{Q})^{-1} \boldsymbol{Z}_S^T)$. The optimal smoothing parameter is chosen to minimize the GCV measure. We refer to the complete estimating procedure as the components selection and smoothing in a sparse functional additive model (CSS-FAM).

Remark: Note that Model 2.3 and the corresponding estimation scheme presented above only account for the effect of a single functional predictor on a scalar response. Examples show that incorporating scalar predictors is likely to improve the prediction accuracy in practice (Sang, Lockhart and Cao (2018); Wong, Li and Zhu (2018)). Therefore, when predicting a scalar response is the main goal, it would be desirable to incorporate both scalar predictors and multiple functional predictors into the current model structure using the adaptive group LASSO and smoothing spline for estimation. Sang, Lockhart and Cao (2018) described how to extend such a framework to allow for scalar covariates and multiple functional predictors.

## 3. Theoretical Properties

To ensure that the estimated scaled FPC scores, $\hat{\boldsymbol{\zeta}}$, are consistent estimators of the true scaled FPC scores, we need to impose regularity conditions on the design of the functional predictor $X(t)$. The following conditions follow Zhu, Yao and Zhang (2014). As stated in Section 2.1, $\{t_{ij}, j = 1, \ldots, N_i; i = 1, \ldots, n\} \subset \mathcal{I}$ denote the time points when the functional predictor $X_i(t)$ is observed. We assume that $t_{i0} = 0$ and $t_{iN_i} = T$ for each $X_i(t)$. Let $\mathcal{I}_\tau = [-\tau, T + \tau]$, for some $\tau > 0$, and let $h_i$ and $K(\cdot)$ denote the bandwidth and the kernel function, respectively, used in smoothing the $i$th trajectory. Note that the same kernel function is employed in the local linear smoother for each trajectory when estimating FPC

scores. The following regularity conditions guarantee that the estimated FPC scores and eigenvalues of the covariance function of $X(t)$ converge in probability to the corresponding population values.

## Assumption A

(A1) $X(t)$ has a continuous second derivative on $\mathcal{I}_d$ with probability 1, and for $k = 0, 2$, $\int \mathrm{E}\{X^{(k)}(t)\}^4 dt < \infty$. The measurement errors $e_{ij}$ of $X_i(t)$ satisfy $\mathrm{E}(e_{ij}^4) < \infty$ and are i.id.

(A2) We define $T_n$ as the lower bound of the number of observations for each trajectory $X_i(t)$. As $n \to \infty$, $T_n \to \infty$. Let $\triangle_i$ denote the largest time difference between two consecutive observations for each trajectory $X_i(t)$; that is, $\triangle_i = \max\{t_{ij} - t_{i,j-1} : j = 1, \ldots, N_i\}$. The maximal value of these satisfies $\max_i \triangle_i = O(T_n^{-1})$.

(A3) There is a sequence $b_n \to 0$, such that $c_1 b_n \leq \min_i h_i \leq \max_i h_i \leq c_2 b_n$, for some constants $c_2 \geq c_1 > 0$, as $n \to \infty$. In addition, $b_n$ and $T_n$ satisfy $(T_n b_n)^{-1} + b_n^4 + T_n^{-2} = O(n^{-1})$.

(A4) The kernel function $K(\cdot)$ has compact support and satisfies $|K(s) - K(t)| \leq C|s - t|$ for $s$, $t$ in its domain and some positive constant $C$.

For Model (2.3), let $A_1$ and $A_0$ denote the set of nonvanishing and vanishing components, respectively; that is, $A_1 = \{j : f_j \neq 0, j = 1, \ldots, d\}$ and $A_0 = \{j : f_j \equiv 0, 1 \leq j \leq d\}$. For the transformation function $F(x|\lambda)$, with a cdf with variance $\lambda$, we make the following assumptions.

## Assumption B

(B1) The transformation function $F(x|\lambda)$ is differentiable at $x$ and $\lambda$. Furthermore, there exist a positive constant $C$ and a negative constant $\gamma$, such that $\partial F(x|\lambda)/\partial x \leq C\lambda^\gamma$ and $\partial F(x|\lambda)/\partial \lambda \leq C\lambda^\gamma|x|$.

(B2) The cdf of each scaled score $\zeta_j$ is absolutely continuous, and there exist positive constants $C_1$ and $C_2$, such that the probability density function of $\zeta_j$, $g_j$, satisfies $C_1 \leq g_j(x) \leq C_2$, for $x \in [0, 1]$ and $j \in A_1$.

Assumption (B1) is from Zhu, Yao and Zhang (2014). Together with Assumptions (A1)-(A4), it guarantees that $\hat{\zeta}_j$ is a consistent estimator of $\zeta_j$, for $1 \leq j \leq d$. Assumption (B2) is a standard assumption in nonparametric additive models, according to Stone (1985).

Define $||f||_2 = \{\int_0^1 f^2(x)dx\}^{1/2}$ whenever the integral is finite. Let $L > 0$, $r$ be a nonnegative integer, and $\nu \in (0,1]$ such that $\rho = r + \nu > 0.5$. Let $\mathscr{F}$ be the class of functions $h$ on $[0,1]$ with the $r$th derivatives that exist and satisfy the Hölder condition with exponent $\nu$: $|h^{(r)}(s) - h^{(r)}(t)| \leq L|s - t|^\nu$, for any $0 \leq s, t \leq 1$. Other standard assumptions for additive nonparametric models (see Huang, Horowitz and Wei (2010)) include:

## Assumption C

(C1) $\min_{j \in A_1} ||f_j|| \geq c_f$, for some $c_f > 0$.

(C2) The random variables $\epsilon_1, \ldots, \epsilon_n$ are i.i.d. with mean zero and variance $\sigma_\epsilon^2$. Furthermore, the tail probability satisfies $P(|\epsilon_1| > x) \leq K \exp(-Cx^2)$, for $\forall x \geq 0$ and some constants $C$ and $K$.

(C3) $\mathrm{E}\{f_j(\zeta_j)\} = 0$ and $f_j \in \mathscr{F}$, for $j \in A_1$.

The following proposition explains why it is reasonable to employ B-spline functions to approximate each nonparametric component $f_j$ in Model (2.3). To guarantee that B-spline functions in $\mathscr{S}_n$ can provide a satisfactory approximation of functions in $\mathscr{F}$, we assume that $l$, the order of the polynomial functions in $\mathscr{S}_n$, satisfies $l > \max\{r, 1\}$. Write the centered version of $\mathscr{S}_n$ as

$$\mathscr{S}_{nj}^0 = \left\{ f_{nj} : f_{nj}(x) = \sum_{k=1}^{m_n} \beta_{jk}\psi_k(x), (\beta_{j1}, \ldots, \beta_{jm_n}) \in \mathbb{R}^{m_n} \right\}, \ j = 1, \ldots, d,$$

where $\psi_k$ is a centered spline basis functions, as defined in Section 2.2.

**Proposition 1.** *Suppose that $f \in \mathscr{F}$ and $\mathrm{E}\,f(\zeta_j) = 0$. Then, under Assumptions A and B, there exists an $f_{nj} \in \mathscr{S}_{nj}^0$, such that*

$$\frac{1}{n}\sum_{i=1}^n \{f_{nj}(\hat{\zeta}_{ij}) - f(\hat{\zeta}_{ij})\}^2 = O_p(m_n^{-2\rho} + n^{-1})$$

*if $m_n = O(n^\alpha)$ with $0 < \alpha < 0.5$.*

Let $\boldsymbol{\psi}(x) = (\psi_1(x), \ldots, \psi_{m_n}(x))^T$ for $x \in [0,1]$. Proposition 1 implies that, uniformly over $j \in \{1, \ldots, d\}$, there exists $\boldsymbol{\beta}_j \in \mathbb{R}^{m_n}$, such that $(1/n)\sum_{i=1}^n \{\boldsymbol{\beta}_j^T \boldsymbol{\psi}(\hat{\zeta}_{ij}) - f(\hat{\zeta}_{ij})\}^2 = O_p(m_n^{-2\rho} + n^{-1})$ under Assumptions A and B, provided $m_n = O(n^\alpha)$. Furthermore, we can take $\boldsymbol{\beta}_j = \boldsymbol{0}$ for $j \in A_0$. Denote $\{j : \tilde{\boldsymbol{\beta}}_j \neq \boldsymbol{0}\}$ and $\{j : \tilde{\boldsymbol{\beta}}_j = \boldsymbol{0}\}$ as $\tilde{A}_1$ and $\tilde{A}_0$, respectively. Theorem 1 establishes the selection

consistency and estimation consistency of $\tilde{\boldsymbol{\beta}}_j$ obtained from the group LASSO step.

**Theorem 1.** *Suppose that Assumptions* A, B, *and* C *hold and* $\lambda_1 \geq C\sqrt{n\log(m_n)}$ *for some sufficiently large constant* $C$. *Then:*

(i) *If* $m_n \to \infty$ *as* $n \to \infty$ *with rate satisfying* $m_n = o(n^{1/6})$, *and* $(\lambda_1^2 m_n^2)/n^2 \to 0$ *as* $n \to \infty$, *then all nonzero* $\boldsymbol{\beta}_j$, *for* $j \in A_1$, *are selected with probability converging to one.*

(ii) *If* $m_n = o(n^{1/6})$, *then* $\sum_{j=1}^{d} ||\tilde{\boldsymbol{\beta}}_j - \boldsymbol{\beta}_j||_2^2 = O_p\left(m_n^2 \log m_n/n\right) + O_p\left(m_n^2 \lambda_1^2/n^2\right) + O_p\left(m_n/n + 1/m_n^{2\rho-1}\right)$.

Theorem 2 further illustrates that the estimated functions obtained from the group LASSO step, $\tilde{f}_j$, enjoy selection consistency and estimation consistency.

**Theorem 2.** *Suppose that Assumptions* A, B, *and* C *hold and* $\lambda_1 \geq C\sqrt{n\log(m_n)}$ *for some sufficiently large constant* $C$. *Then:*

(i) *If* $m_n \to \infty$ *as* $n \to \infty$ *with rate satisfying* $m_n = o(n^{1/6})$, *and* $(\lambda_1^2 m_n)/n^2 \to 0$ *as* $n \to \infty$, *then in the group LASSO step, all the nonzero additive components* $f_j$, *for* $j \in A_1$, *are selected with probability converging to one.*

(ii) *If* $m_n = o(n^{1/6})$, *then* $||\tilde{f}_j - f_j||_2^2 = O_p\left(m_n \log m_n/n\right) + O_p\left(m_n/n + 1/m_n^{2\rho}\right) + O_p\left(m_n\lambda_1^2/n^2 + m_n/n\right)$, *for* $j \in A_1 \cup \tilde{A}_1$.

For two (positive) sequences $\{a_n\}$ and $\{b_n\}$, if $a_n/b_n$ is bounded away from 0 and $\infty$, then denote this as $a_n \sim b_n$. The following corollary is derived directly from Theorem 2.

**Corollary 1.** *Suppose that Assumptions* A, B, *and* C *hold. If* $m_n \sim n^{1/(2\rho+1)}$ *and* $\lambda_1 \sim \sqrt{n\log(m_n)}$, *then:*

(i) *If* $\rho > 5/2$, *then in the group LASSO step, all nonzero additive components* $f_j$, *for* $j \in A_1$, *are selected with probability converging to one.*

(ii) *If* $\rho > 5/2$, *then* $||\tilde{f}_j - f_j||_2^2 = O_p(n^{-2\rho/(2\rho+1)} \log m_n)$, *for* $j \in A_1 \cup \tilde{A}_1$.

Theorem 3 states that the adaptive group LASSO yields an estimate that is also consistent in both selection and estimation. Furthermore, it illustrates that this estimate compares favorably with that given by the group LASSO with respect to estimation accuracy.

**Theorem 3.** *Suppose that Assumptions* A, B, *and* C *hold and* $m_n \sim n^{1/(2\rho+1)}$, *where* $\rho > 5/2$. *If the tuning parameters satisfy* $\lambda_1 \sim \sqrt{n \log(m_n)}$, $\lambda_2 \leq O(n^{\frac{1}{2}})$, $\lambda_2/n^{(8\rho+3)/(8\rho+4)} = o(1)$, *and* $n^{1/(4\rho+2)}\sqrt{\log(m_n)}/\lambda_2 = o(1)$, *then we have the following*

   (i) *With probability approaching one, the nonzero components* $\{f_j, j \in A_1\}$ *are selected and* $||\hat{f}_j||_2 = 0$, *for* $j \in A_0$.

  (ii) $\sum_{j \in A_1} ||\hat{f}_j - f_j||_2^2 = O_p(n^{-2\rho/(2\rho+1)})$.

## 4. Simulation Studies

In this section, we use simulated examples to illustrate the properties of our proposed estimator. We also compare our method with several conventional methods commonly used in practice.

We simulate the data as follows. In each simulation replicate, we generate $n$ curves, and the observations are made at $m = 200$ equally spaced points in $[0, 10]$. In our simulation studies, we set $n = 100$ or $500$. To accommodate measurement errors, the observation at $t_j$ $(j = 1, \ldots, m)$ is generated as $W_{ij} = X_i(t_j) + e_{ij}$, where $\{X_i(t)\}_{i=1}^n$ are i.i.d. samples of a stochastic process $X(t)$, and $e_{ij}$ are i.i.d. normals with mean zero and variance $0.1$. For $k = 1, \ldots, 20$, let $\lambda_k = 31.5 \times 0.6^k$ denote the $k$th eigenvalue of the covariance function of $X(t)$. The corresponding $k$th eigenfunction is the $k$th Fourier basis function, denoted by $\phi_k(t)$. Then, $X_i(t) = m(t) + \sum_{k=1}^{20} \xi_{ik}\phi_k(t)$, where $m(t) = t + \sin t$ denotes the mean function of $X(t)$, and $\{\xi_{ik}\}_{k=1}^{20}$ are independently sampled from $N(0, \lambda_k)$. The scaled score $\zeta_{ik}$ is defined as the uniform score of $\xi_{ik}$; that is, $\zeta_{ik} = \Phi(\xi_{ik}/\sqrt{\lambda_k})$, for $k = 1, \ldots, 20$, $i = 1, \ldots, n$, where $\Phi$ denotes the cdf of a standard normal distribution. The response variable is generated from Model (2.3): $y_i = a + f_1(\zeta_{i1}) + f_2(\zeta_{i2}) + f_4(\zeta_{i4}) + \epsilon_i$. We set the true intercept to $a = 1.2$, and the true components to $f_1(x) = x\exp(x) - 1$, $f_2(x) = \cos(2\pi x)$ and $f_4(x) = 3(x - (1/4))^2 - 7/16$, for $x \in [0, 1]$. The random errors $\epsilon_i$ are independently sampled from a normal distribution with mean zero and variance $0.67$. The signal-to-noise ratio is defined as $\mathrm{Var}\{f_1(\zeta_1) + f_2(\zeta_2) + f_4(\zeta_4)\}/\mathrm{Var}(\epsilon)$; we set this ratio to be approximately two. We estimate the model by fitting $n$ randomly generated training observations, and evaluate its performance on $200$ randomly generated test observations. The simulation is implemented for $100$ simulation replicates. The simulation results for $n = 200$ and $300$ are presented in the Supplementary Material.

In addition to employing the proposed method CSS-FAM, we fit the data us-

ing three conventional models: multivariate adaptive regression splines (MARS) (Friedman, Hastie and Tibshirani (2001)), two extended functional additive models (FAM) proposed by Müller and Yao (2008), and the component selection and estimation for the functional additive model (CSE-FAM) proposed by Zhu, Yao and Zhang (2014). More specifically, MARS is fitted using the function **earth** in the R package **earth**, and the variables that are included in the final model are examined by the function **evimp**. In the first extended FAM, denoted by FAM, the response variable $y$ is fitted with a multiple linear regression, where the covariates are $f_1(\hat{\zeta}_1)$, $f_2(\hat{\zeta}_2)$, and $f_4(\hat{\zeta}_4)$. In other words, FAM assumes to know the true model structure based on three true covariates $f_1(\hat{\zeta}_1)$, $f_2(\hat{\zeta}_2)$, and $f_4(\hat{\zeta}_4)$. The second extended FAM, denoted by S-FAM, uses a saturated model to incorporate the first $d$ FPC scores, such that they explain over 99.9% of the total variability in the smoothed sample curves. The value of $d$ is 15, 16, or 17 in all simulation replicates. We employ the function **gam** in the R package **mgcv** to fit this model, in which the arguments of the additive components are $\hat{\zeta}_j$, for $j = 1, \ldots, d$. Then p-values of all terms in the model are available from the function **summary.gam**. Only the significant nonparametric components (p-value $< 0.05$) are retained when computing the true positive (TP) rate and the false positive (FP) rate. We also consider an alternative method for estimating Model (2.3), a group LASSO and an adaptive group LASSO, denoted by AGL-FAM.

Table 1 summarizes the comparison between these six methods based on $1,000$ simulation replicates. The results suggest that, compared with CSE-FAM, CSS-FAM has similar performance in prediction when the sample size $n = 100$ or 500. Both outperform the other three methods, except FAM, in terms of prediction accuracy, and are slightly inferior to FAM, which assumes the true components are known. This suggests that the extra adaptive smoothing spline step can increase the prediction accuracy when an adaptive group LASSO yields a wiggly estimate. In terms of the quality of the estimations of the nonparametric components, CSS-FAM can rival CSE-FAM as well, because both yield estimates that are reasonably close to the true nonparametric components. In addition, the residual sum of squares (RSS) for each component estimated using CSS-FAM is much smaller than that using AGL-FAM, indicating that the smoothing spline enables us to obtain a smoother and more accurate estimate of the nonparametric components.

Table 1 also compares these methods from the perspective of variable selection, where the true positive (TP) rate and the false positive (FP) rate are employed for assessment. Combined with Table S1 in the Supplementary Ma-

Table 1. Summary statistics for the six methods. MSPE refers to the mean squared prediction error on the test data; the residual sum of squares (RSS) for each estimated component $\hat{f}_j$ is defined as: $\mathrm{RSS}(\hat{f}_j) = \int_0^1 (\hat{f}_j(x) - f_j(x))^2 dx$; TP% and FP% stand for the true positive and false positive rates, as a percentage, respectively. The point estimate for each measure is averaged over 100 simulation replicates, and the corresponding estimated standard error is given in parentheses.

| Statistics | $n$ | Methods | | | | | |
|---|---|---|---|---|---|---|---|
| | | MARS | FAM | S-FAM | CSE-FAM | AGL-FAM | CSS-FAM |
| MSPE | 100 | 1.15 (0.25) | 0.92 ( 0.17) | 1.15 ( 0.21) | 1.00 ( 0.19) | 1.25 (0.23 ) | 1.01 ( 0.20 ) |
| | 500 | 0.78 (0.09) | 0.73 ( 0.08) | 0.77 ( 0.09) | 0.74 ( 0.08) | 0.80 (0.09 ) | 0.74 ( 0.08 ) |
| $\mathrm{RSS}(\hat{f}_1)$ | 100 | - | 2.6 ( 4.9 ) | 3.6 ( 4.7 ) | 2.6 ( 4.1 ) | 13.4 (6.9 ) | 3.8 ( 5.7 ) |
| $(\times 10^{-2})$ | 500 | - | 0.4 ( 0.5 ) | 0.6 ( 0.6 ) | 0.5 ( 0.4 ) | 2.5 (0.9 ) | 0.6 ( 0.4 ) |
| $\mathrm{RSS}(\hat{f}_2)$ | 100 | - | 6.8 (10.5 ) | 11.2 (10.0 ) | 18.1 (13.6 ) | 12.8 (6.8 ) | 8.1 (13.8 ) |
| $(\times 10^{-2})$ | 500 | - | 0.5 ( 0.7 ) | 1.9 ( 1.3 ) | 2.9 ( 1.5 ) | 3.1 (1.4 ) | 1.9 ( 1.3 ) |
| $\mathrm{RSS}(\hat{f}_4)$ | 100 | - | 6.7 (10.3 ) | 4.0 ( 3.3 ) | 4.6 ( 5.7 ) | 14.3 (7.2 ) | 5.9 (11.2 ) |
| $(\times 10^{-2})$ | 500 | - | 0.7 ( 1.1 ) | 0.7 ( 0.5 ) | 0.5 ( 0.4 ) | 2.3 (1.0 ) | 0.5 ( 0.7 ) |
| TP% | 100 | 99.1 (0.05) | - | 98.2 ( 0.08) | 95.7 ( 0.12) | 94.7 (0.17 ) | 94.7 ( 0.10 ) |
| | 500 | 100 (0.00) | - | 100 ( 0.0 ) | 100 ( 0.0 ) | 100 (0.0 ) | 100 ( 0.0 ) |
| FP% | 100 | 20.4 (0.12) | - | 13.7 ( 0.11) | 3.8 ( 0.07) | 0.9 (0.03 ) | 0.9 ( 0.03 ) |
| | 500 | 29.0 (0.14) | - | 8.9 ( 0.08) | 3.0 ( 0.07) | < 0.01 (0.003) | <0.01 ( 0.003) |
| Time | 100 | 0.01 (0.03) | < 0.01 | 0.39 ( 0.21) | 2.87 ( 0.23) | 0.48 (0.04 ) | 2.40 ( 0.10 ) |
| (seconds) | 500 | 0.02 (0.06) | < 0.01 | 2.88 ( 2.28) | 117.2 ( 5.91) | 3.57 (0.27 ) | 11.4 ( 2.77 ) |

terial, we find that although CSS-FAM and AGL-FAM perform slightly worse than the other models in terms of correctly selecting nonzero variables for a relatively small sample ($n = 100$ or $200$), this tiny gap vanishes when the sample size increases ($n = 300$ or larger). Furthermore, CSS-FAM and AGL-FAM dominate the other methods in not selecting irrelevant components, regardless of how large or small the sample size is. The other methods mistakenly select irrelevant variables far more often than CSS-FAM or AGL-FAM does, especially when the sample size is relatively small.

The computational time for each method is recorded in Table 1 as well. Obviously, CSE-FAM is the most computationally intensive method if a full basis is employed. This is a serious issue in implementations, particularly when the sample size is large, as mentioned in Section 1. In comparison, the proposed method, CSS-FAM, can still be implemented within 12 seconds, even when the training data set consists of 500 curves.

Figure S1 in the Supplementary Material illustrates the estimation details

for one randomly selected simulation replicate when the number of curves is $n = 500$. After estimating the scaled FPC score, we fit a group LASSO on the training data, as shown in (2.5). The top, left panel in Figure S1 describes how the five-fold cross-validation error changes with $\lambda_1$. The optimal $\lambda_1$ is chosen to minimize this. The middle panel explains how to choose the optimal $\lambda_2$ for the adaptive group LASSO step in (2.6), based on five-fold cross validation. The top, right panel shows how to choose the optimal smoothing parameter ($\lambda_3$) by minimizing GCV in the smoothing spline step. The bottom three panels in Figure S1 illustrate the effects of the extra smoothing spline step on the estimation of the nonparametric components after using adaptive group LASSO. The adaptive group LASSO method may lead to an excessively wiggly estimate for each nonzero nonparametric component. Smoothing splines can control this roughness and, hence, yield a smoother and more accurate estimate.

## 5. Applications

In this section, we fit the sparse functional additive model (2.3) using our proposed method (CSS-FAM), together with several conventional models considered in the simulation studies, to analyze two real data sets. An application to air pollution data is introduced in the Supplementary Material. In addition to the models considered in the simulation, we fit a multiple linear model to investigate whether a functional linear model can adequately characterize the relationship between the scalar response and the functional predictor in these two examples. The covariates in the multiple linear model are the first $d$ FPC scores. We choose the truncation level $d$ in the same way as for Model (2.2). This multiple linear model is actually a special case of Model (2.2): each additive component takes a linear form. The LASSO (Tibshirani (1996)) is implemented when fitting the mulitple linear model in these two examples to obtain a more parsimonious model and to reduce the variability. We refer to this method as the LAF. In the air pollution data, the trajectories of the functional predictor for some subjects are sparsely observed. In contrast, in the Tecator data, the functional predictor is regularly spaced and densely observed across all subjects. In each example, we randomly divide the whole data set into a training set and a test set. The training set is used to fit each model, and the test set is used for evaluation. All models are compared with respect to the mean squared prediction error.

Table 2. Mean squared prediction errors (MSPEs) on the test data for six methods.

| Methods | MARS | LAF | S-FAM | CSE-FAM | AGL-FAM | CSS-FAM |
|---------|------|-----|-------|---------|---------|---------|
| MSPE    | 0.99 | 0.66 | 0.56 | 0.55   | 0.92    | 0.51    |

## 5.1. Tecator data

The Tecator data are recorded for 240 meat samples on a Tecator Infratec Food and Feed Analyzer working in the wavelength range 850-1050 nm using the near infrared transmission (NIT) principle. Each record contains of a 100-channel spectrum of absorbance, and the percentages of three components of the meat: moisture (water), fat, and protein. The spectrum records the negative base 10 logarithm of the transmittance measured by the spectrometer. The percentages of the three meat components are determined by analytic chemistry. As demonstrated by a large body of research (see Vila, Wagner and Neveu (2000), Goldsmith and Scheipl (2014), Zhu, Yao and Zhang (2014)), the spectrum of absorbance is highly predictive of the percentage of these three meat components. Figure S2 in the Supplementary Material shows the trajectories of the spectrum of absorbance of the 240 meat samples. Here, we examine the effect of the spectral trajectories of the meat sample on the protein content using the sparse functional additive model (2.3).

The protein content, denoted by $Y$, is the response variable of primary interest; the functional predictor $X(t)$ denotes the spectrum of absorbance. An FPCA is performed to estimate the FPC scores, and then to obtain the scaled FPC scores, denoted by $\hat{\zeta} = (\hat{\zeta}_1, \ldots, \hat{\zeta}_d)$. Zhu, Yao and Zhang (2014) suggested that the first $d = 20$ should be retained in order to achieve satisfactory prediction accuracy, even though the first 10 FPCs explain more than 99.9% of the total variability in the smoothed sample curves. To compare the performance of the various methods in terms of their prediction accuracy, the 240 meat samples are divided into a training sample and a test sample. According to the original description of the data set (`http://lib.stat.cmu.edu/datasets/tecator`), the 240 meat samples are divided into three categories: a training set of 172 meat samples, a test set of 43 meat samples, and 25 samples that are used for extrapolation, and should be ignored. We, however, randomly choose 187 meat samples to train the model; the test set comprises the remaining 53 meat samples.

The six models are compared with respect to their prediction accuracy in Table 2. Clearly, CSS-FAM outperforms the other methods in terms of prediction. In particular, the difference between CSS-FAM and LAF implies that a

linear model cannot adequately characterize the relationship between the protein content and the spectrum of absorbance of the meat samples. Nevertheless, CSS-FAM achieves a better tradeoff between flexibility and simplicity than other methods do. Additionally, the poor performance of AGL-CSS, especially when compared with CSS-FAM, suggests that the extra smoothing spline step in the proposed algorithm enhances the prediction accuracy considerably.

In AGL-FAM, 10 cubic B-spline basis functions are employed to represent the nonparametric components in the sparse functional additive model (2.3). A five-fold cross-validation suggests that $\lambda_1 = 0.002$ is an optimal choice of the penalty parameter in the group LASSO step, and $\lambda_2 = 0.011$ minimizes the five-fold cross-validation error in the adaptive LASSO step. As a result, 14 nonvanishing components, $\{\hat{f}_1, \ldots, \hat{f}_9, \hat{f}_{11}, \hat{f}_{16}, \hat{f}_{17}, \hat{f}_{19}, \hat{f}_{20}\}$, are selected from the 20 components. This finding is slightly inconsistent with the conclusion drawn in Zhu, Yao and Zhang (2014), who claim that $\{\hat{f}_1, \ldots, \hat{f}_8, \hat{f}_{10}, \hat{f}_{13}, \hat{f}_{16}, \hat{f}_{17}\}$ are nonvanishing components. To refine these estimated components, a smoothing spline is employed and the optimal choice of smoothing parameter, $\lambda_3 = 0.001$, is chosen to minimize the GCV measure.

## 6. Conclusion

Compared with the traditional FLR, the sparse functional additive model (2.3) proposed in this article provides a more flexible description of the relationship between a scalar response and a functional predictor. To achieve sparseness, we employ the group LASSO penalty to select and estimate nonzero components in the nonparametric additive model, thereby reducing variability and enhancing interpretability.

The estimation procedure consists of several important techniques. An FPCA is employed to estimate the FPC scores and eigenvalues of the covariance function of the functional predictor. Then, we use B-spline basis functions to represent the nonparametric additive components in the sparse functional additive model (2.3). The use of the group LASSO penalty enables us to select and estimate the nonzero components. To obtain a better estimate of the coefficient vectors, we use the adaptive group LASSO to allow the shrinkages to vary by component. Note that the estimated components given by the adaptive LASSO may not be smooth, because a large number of B-spline basis functions are used to represent them. Thus, we propose using smoothing splines to further refine the estimated nonzero components obtained from the group LASSO step. Simulation

studies demonstrate that this smoothing step improves both the estimation of the additive components and the prediction of the response.

We justify theoretically that our proposed estimator enjoys both selection consistency and estimation consistency. These consistency results are also demonstrated by simulation studies. Two real-data applications show that the proposed model, together with the estimating method, provides an appealing tool for predicting a scalar response from a functional predictor.

Even though we regress a scalar response on a functional covariate only, the methodology can be extended to accommodate other scenarios. For example, this framework can be extended to explore the relationship between a scalar response, whose distribution belongs to the exponential family, and a functional predictor. In addition, in this work, the truncation level $d$, such that the first $d$ FPCs explain over 99.9% of total variability in the functional predictor, is assumed to be fixed. From a theoretical perspective, it is worthwhile investigating the properties of the corresponding estimator when $d$ is allowed to increase with the sample size; this is left to future work.

## Supplementary Material

The online Supplementary Materials describes the procedure used to estimate FPC scores in Section 2.1, and provides proofs of the theoretical results in Section 3. Additional simulation results appear in Section 4 and, in Section 5 we present an example in which we apply the proposed model. The R code for our real-data analysis and the simulation studies can be downloaded at `https://github.com/caojiguo/fam`.

## Acknowledgments

## References

Ainsworth, L. M., Routledge, R. and Cao, J. (2011). Functional data analysis in ecosystem research: the decline of oweekeno lake sockeye salmon and wannock river flow. *Journal of Agricultural, Biological, and Environmental Statistics* **16**, 282–300.

Chen, D., Hall, P., Müller, H.-G. et al. (2011). Single and multiple index functional regression models with nonparametric link. *The Annals of Statistics* **39**, 1720–1747.

De Boor, C. (2001). *A Practical Guide to Splines*. Springer-Verlag, New York.

Eilers, P. H. and Marx, B. D. (1996). Flexible smoothing with B-splines and penalties. *Statistical Science* **11**, 89–102.

Friedman, J., Hastie, T. and Tibshirani, R. (2001). *The Elements of Statistical Learning*. Springer, Berlin.

Goldsmith, J. and Scheipl, F. (2014). Estimator selection and combination in scalar-on-function regression. *Computational Statistics & Data Analysis* **70**, 362–372.

Hastie, T. and Tibshirani, R. (1986). Generalized additive models. *Statistical Science* **1**, 297–310.

Horváth, L. and Kokoszka, P. (2012). *Inference for Functional Data with Applications*. Springer, New York.

Huang, J., Horowitz, J. L. and Wei, F. (2010). Variable selection in nonparametric additive models. *The Annals of Statistics* **38**, 2282.

Lin, Y. and Zhang, H. H. (2006). Component selection and smoothing in multivariate nonparametric regression. *The Annals of Statistics* **34**, 2272–2297.

Lin, Z., Cao, J., Wang, L. and Wang, H. (2017). Locally sparse estimator for functional linear regression models. *Journal of Computational and Graphical Statistics* **26**, 306–318.

Liu, B., Wang, L. and Cao, J. (2017). Estimating functional linear mixed-effects regression models. *Computational Statistics & Data Analysis* **106**, 153–164.

Luo, W., Cao, J., Gallagher, M. and Wiles, J. (2013). Estimating the intensity of ward admission and its effect on emergency department access block. *Statistics in Medicine* **32**, 2681–2694.

Morris, J. S. (2015). Functional regression. *Annual Review of Statistics and Its Application* **2**, 321–359.

Müller, H.-G. and Stadtmüller, U. (2005). Generalized functional linear models. *The Annals of Statistics* **33**, 774–805.

Müller, H.-G., Wu, Y. and Yao, F. (2013). Continuously additive models for nonlinear functional regression. *Biometrika* **100**, 607–622.

Müller, H.-G. and Yao, F. (2008). Functional additive models. *Journal of the American Statistical Association* **103**, 1534–1544.

Ramsay, J. O. and Silverman, B. W. (2005). *Functional Data Analysis*, 2nd Edition. Springer-Verlag, New York.

Reiss, P. T., Goldsmith, J., Shang, H. L. and Ogden, R. T. (2017). Methods for scalar-on-function regression. *International Statistical Review* **85**, 228–249.

Sang, P., Lockhart, R. A. and Cao, J. (2018). Sparse estimation for functional semiparametric additive models. *Journal of Multivariate Analysis* **168**, 105–118.

Stone, C. J. (1985). Additive regression and other nonparametric models. *The Annals of Statistics* **13**, 689–705.

Stone, C. J. (1986). The dimensionality reduction principle for generalized additive models. *The Annals of Statistics* **14**, 590–606.

Tibshirani, R. (1996). Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society. Series B (Statistical Methodological)* **58**, 267–288.

Vila, J.-P., Wagner, V. and Neveu, P. (2000). Bayesian nonlinear model selection and neural networks: a conjugate prior approach. *IEEE Transactions on Neural Networks* **11**, 265–278.

Wong, R. K., Li, Y. and Zhu, Z. (2018). Partially linear functional additive models for multivariate functional data. *Journal of the American Statistical Association.* DOI: 10.1080/01621459.2017.1411268.

Wu, H., Lu, T., Xue, H. and Liang, H. (2014). Sparse additive ordinary differential equations for dynamic gene regulatory network modeling. *Journal of the American Statistical Association* **109**, 700–716.

Yao, F. and Müller, H.-G. (2010). Functional quadratic regression. *Biometrika* **97**, 49–64.

Yuan, M. and Lin, Y. (2006). Model selection and estimation in regression with grouped variables. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* **68**, 49–67.

Zhang, H. H. and Lin, Y. (2006). Component selection and smoothing for nonparametric regression in exponential families. *Statistica Sinica* **16**, 1021–1041.

Zhu, H., Yao, F. and Zhang, H. H. (2014). Structured functional additive regression in reproducing kernel Hilbert spaces. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* **76**, 581–603.

Zou, H. (2006). The adaptive LASSO and its oracle properties. *Journal of the American statistical association* **101**, 1418–1429.

Department of Statistics and Actuarial Science, University of Waterloo, 200 University Avenue West, Waterloo, ON N2L 3G1, Canada.

E-mail: psang@uwaterloo.ca

Department of Statistics and Actuarial Science, Simon Fraser University, Burnaby, BC, V5A 1S6, Canada.

E-mail: liangliang_Wang@sfu.ca

Department of Statistics and Actuarial Science, Simon Fraser University, Burnaby, BC, V5A 1S6, Canada.

E-mail: jiguo_cao@sfu.ca