

Statistica Sinica Preprint No: SS-2018-0483

Title	Robust Inference in Varying-coefficient Additive Models for Longitudinal/Functional Data
Manuscript ID	SS-2018-0483
URL	http://www.stat.sinica.edu.tw/statistica/
DOI	10.5705/ss.202018.0483
Complete List of Authors	Lixia Hu Tao Huang and Jinhong You
Corresponding Author	Tao Huang
E-mail	huang.tao@mail.shufe.edu.cn
Notice: Accepted version subject to English editing.	

ROBUST INFERENCE IN VARYING-COEFFICIENT ADDITIVE MODELS FOR LONGITUDINAL/FUNCTIONAL DATA

Lixia Hu¹, Tao Huang² and Jinhong You²

2

¹ *Shanghai Lixin University of Accounting and Finance*

3

² *Shanghai University of Finance and Economics*

4

Abstract: The focus in this paper is on the robust inference of

5

a varying-coefficient additive model for sparse or dense longitudi-

6

nal/functional data. A spline-based three-step M-estimation method

7

is proposed for estimating varying-coefficient component function-

8

s and additive component functions. Above all, the consistency

9

and asymptotic normality of sparse data and dense data are stud-

10

ied in a unified framework. Furthermore, employing a regularized

11

M-estimation method, a model identification procedure is proposed

12

that consistently identifies an additive term and a varying-coefficient

13

term. Simulation studies investigating the finite-sample performance

14

of the proposed methods confirm the asymptotic theory, and real-life

15

examples illustrating the proposed methods are considered.

16 *Key words and phrases:* B-spline, M-estimator, SCAD, tensor prod-
17 uct, varying-coefficient additive model.

18 **1. Introduction**

19 Repeated-measurement data arise often in the clinical, biometrical, epidemi-
20 ological, social, and economic fields (Diggle, Liang and Zeger , 1994). Common
21 among such data are longitudinal and functional data, which have different sam-
22 pling mechanisms. Typically, data are termed functional when they are recorded
23 densely over time in a continuum without noise, and they are termed longitudinal
24 when the measurements are made at a few discrete time points with experimen-
25 tal error. However, in practice functional data are analyzed after smoothing
26 noisy observations (Ramsay and Ramsey , 2002). A vast body of existing liter-
27 ature considers the statistical inference of functional data based on observations
28 at discrete time points contaminated with measurement errors, a practice that
29 makes it possible to analyze longitudinal data and functional data in a unified
30 framework (Li and Hsing, 2010; Yao , 2007). Some scholars have also studied
31 longitudinal data by means of functional principal components analysis (Yao,
32 Müller and Wang , 2005).

33 In a typical repeated-measurement-data setting, a sample of n subjects or
34 curves is observed at n_i discrete time points. If each n_i exceeds some power
35 of n , then the data are referred to as dense data. If each n_i is bounded by a

36 finite positive number or follows a fixed distribution, then the data are referred
37 to as sparse data. Recently, Zhang and Wang (2016) considered nonparametric
38 estimation of the mean function and covariance function for sparse and dense
39 functional data in a unified framework, and they partitioned the data into sparse,
40 dense, and ultra-dense according to the magnitude of n_i relative to n .

41 There is much literature on nonparametric regression methods for functional
42 data and longitudinal data with sparsity or/and denseness. Because of their sim-
43 plicity, flexibility, and interpretability, varying-coefficient models (VCMs) have
44 been used extensively to analyze longitudinal data (Hoover et al. , 1998; Xue and
45 Zhu , 2007). Additive models (AMs) provide an alternative regression method
46 (Carroll et al. , 2009; Xue, Qu and Zhou , 2010), and Zhang, Park and Wang
47 (2013) proposed a time-varying additive model for analyzing longitudinal da-
48 ta to capture dynamic effects. Recently, for analyzing functional data, Zhang
49 and Wang (2015) proposed a novel nonparametric regression method called
50 the varying-coefficient additive model (VCAM), which covers classical AMs and
51 VCMs as its special cases. Specifically, let $Y(t)$ be a smooth random response
52 process and $\mathbf{X} = (X_1, \dots, X_p)^\tau$ be the p -vector of covariates. The regression
53 function $m(t, \mathbf{x}) := E[Y(t)|\mathbf{X} = \mathbf{x}]$ of a VCAM has the form

$$m(t, \mathbf{x}) = \alpha_0(t) + \sum_{k=1}^p \alpha_k(t)\beta_k(x_k), \quad (1.1)$$

54 where α_k is the varying-coefficient component function and β_k is the additive

55 component function.

56 Zhang and Wang (2015) proposed a two-step spline estimation method for
57 varying-coefficient component functions and additive component functions based
58 upon two key assumptions, namely (i) each subject (smooth process or function
59 curve) is observed at dense time points and (ii) each predictor is subject specific
60 but independent of the observation time. The above conditions are easily satisfied
61 for functional data but are restrictive for longitudinal data. Furthermore, if
62 conditions (i) and/or (ii) are violated, then the estimation method of Zhang and
63 Wang either fails or performs poorly, as shown in Table 6 of the supplementary
64 material. In the present paper, we consider two real-life examples, namely the
65 CD4 cell count in HIV seroconversion (Zeger and Diggle, 1994) and the cigarette
66 dataset from the R package “phtt” (Bada and Liebl, 2012) investigated in the
67 supplementary material. Note that each example violates condition (i) and/or
68 (ii), meaning that the two-step spline estimation method proposed by Zhang and
69 Wang (2015) is not appropriate. One of our aims herein is to relax conditions
70 (i) and (ii) and develop a general estimation method that has wider applications
71 in practical fields.

72 Much of the literature is focused on the classical mean regression method,
73 but that method is sensitive to outliers and non-normal error distributions. An
74 alternative is the M-type robust regression method, which can treat mean, me-

75 dian, quantile, and more-general robust-type regression methods in a unified
76 framework. Many scholars have considered robust regression techniques, such as
77 Koenker and Bassett (1978) for quantile regression of linear models, He and Shi
78 (1994) and He, Zhu and Fung (2002) for M-estimators of partially linear models,
79 and Tang and Cheng (2008) for M-estimators of varying-coefficient models.

80 Herein, we consider the robust inference of a VCAM for sparse and dense
81 longitudinal or functional data, allowing the predictors to be smooth process-
82 es covering condition (ii). We propose spline-based three-step M-estimators for
83 varying-coefficient component functions and additive component functions. The
84 asymptotic properties of the newly-proposed estimators are presented in a unified
85 framework, and we can separate sparse data and dense data according to the rela-
86 tive order of n_i to n , which can be viewed as a generalization of Zhang and Wang
87 (2016) to a VCAM. Similar to Hu, Huang and You (2018), a remarkable aspect
88 of our estimators is the oracle property, which implies that the iteration step
89 does not cause additional asymptotic errors. Furthermore, from the perspective
90 of model parsimony, we develop a spline-based penalized M-estimator to decide
91 whether the product term in (1.1) reduces to a varying-coefficient term or an
92 additive term, corresponding to an additive component function of linear form
93 or a constant varying-coefficient component function, respectively. We also show
94 that an additive term and a varying-coefficient term can be selected correctly

95 with probability approaching unity under mild conditions.

96 The remainder of this paper is organized as follows. In Section 2, we de-
97 scribe the model setup and propose the spline-based three-step M-estimators of
98 univariate component functions. In Section 3, we present the asymptotic the-
99 ories of the proposed estimators. In Section 4, we introduce a robust model
100 identification procedure, and in Section 5 we address the selection of the smooth-
101 ing parameters. In Section 6, we consider simulation examples investigating the
102 finite-sample performance and empirical examples illustrating our method. Fi-
103 nally, in Section 7 we make some brief concluding remarks. The technical proofs
104 and additional numerical studies are included in the supplementary material.

105 **2. Model and Estimation Method**

106 **2.1. Model assumptions**

107 Let $Y(t)$ be a smooth response process and $\mathbf{X}(t) = \{X_1(t), \dots, X_p(t)\}^\tau$ be
108 a p -vector of smooth processes of covariates, where the superscript τ denotes
109 the transpose of a vector or matrix. Without loss of generality, we assume that
110 the response and covariates from a subject are L_2 integrable stochastic processes
111 on the interval $[0, 1]$. The relationship between the response and covariates is
112 modeled by a VCAM as follows:

$$Y(t) = \alpha_0(t) + \sum_{k=1}^p \alpha_k(t) \beta_k(X_k(t)) + U(t), \quad (2.1)$$

113 where $U(t)$ is the stochastic component of response process $Y(t)$, independent of

114 covariate process $\mathbf{X}(t)$, with mean function $E[U(t)] = 0$ and auto-covariance func-
115 tion $\gamma(t, s) = E[U(t)U(s)]$. To uniquely identify univariate component functions,
116 we impose the identification conditions $\int_0^1 \alpha_k(t)dt = 1$ and $E[\beta_k(X_k(t))] = 0$
117 following the common practice in nonparametric regression (Zhang and Wang ,
118 2015; Wang and Yang , 2007; Vogt , 2012; Hu, Huang and You , 2018).

119 In practical applications, the process Y is not observable but can be measured
120 at any given time point with random error e such that $E(e) = 0$, $\text{Var}(e) = \sigma_e^2$.
121 We sample n subjects independently and observe subject i at n_i time points
122 $(t_{i1}, \dots, t_{in_i})$, denoting y_{ij} and $\mathbf{x}_{ij} = (x_{ij1}, \dots, x_{ijp})^\tau$ as the observations of the
123 response and the vector of covariates at time t_{ij} , respectively. Then, the sample
124 version of VCAM (2.1) can be written as

$$y_{ij} = \alpha_0(t_{ij}) + \sum_{k=1}^p \alpha_k(t_{ij})\beta_k(x_{ijk}) + U_{ij} + e_{ij}, \quad (2.2)$$

125 where $U_{ij} = U_i(t_{ij})$ is a realization of the subject-specific random trajectory
126 $U_i(t)$ at observation time t_{ij} , and e_{ij} are independent and identical copies of
127 the random measurement error e . As done by Zhang and Wang (2015), we
128 ignore the intra-subject covariance structure and incorporate the covariance of
129 $\{U_{ij}, j = 1, \dots, n_i\}$ into the random error term, denoted as $\varepsilon_{ij} = U_{ij} + e_{ij}$.

130 **Remark 1.** The product term $\alpha_k(t)\beta_k(x_k)$ in VCAM (2.1) reduces to an additive
131 term if α_k is a constant and to a varying-coefficient term if β_k is a linear function.

132 In other words, a VCAM is more flexible than either an AM or a VCM and can
133 greatly reduce the systematic bias of modelling.

134 2.2. Three-step M-estimation method

135 The spline method is a useful tool for fitting smooth nonparametric func-
136 tions, and the B-spline basis is preferred for its computational stability. Let
137 $\{\tilde{b}_1(x), \dots, \tilde{b}_{K+m}(x)\}$ be a normalized m -order B-spline basis with K interior knots
138 (De Boor, 1978). The scaled version of $\tilde{b}_k(x)$ is given by $b_k(x) = \sqrt{K+m}\tilde{b}_k(x)$,
139 whose favorable properties are presented in the supplementary material. Fur-
140 thermore, similar to Wang and Yang (2007), we construct the centralized ver-
141 sion represented as $\{B_1(x), \dots, B_{K+m-1}(x)\}$. Under the assumption that both
142 $\alpha_k(\cdot)$ and $\beta_k(\cdot)$ are $r(\leq q)$ -order smooth, we adopt a q -order B-spline function
143 to fit a univariate nonparametric function. For any $t \in [0, 1]$ and x in the do-
144 main of $\beta_k(\cdot)$, we use the B-spline bases $\mathbf{b}_C(t) = \{b_1(t), \dots, b_{J_C}(t)\}^\tau$ to approx-
145 imate the varying-coefficient component function $\alpha_k(t)$, and we use $\mathbf{B}_{k,A}(x) =$
146 $\{B_{k1}(x), \dots, B_{kJ_A}(x)\}^\tau$ to approximate the additive component function $\beta_k(x)$ for
147 each $k = 1, \dots, p$, where J_C and J_A are a sum of smooth degree r and the number
148 of interior knots, respectively. The tensor product of $\mathbf{B}_{k,A}(x_k)$ and $\mathbf{b}_C(t)$ is de-
149 fined as $\mathcal{T}_k(t, x_k) = \mathbf{B}_{k,A}(x_k) \otimes \mathbf{b}_C(t)$, where \otimes represents the Kronecker product
150 of matrices or vectors.

151 Now, we propose a spline-based three-step M-estimation method. Specifi-

152 cally, we begin by obtaining initial estimators of varying-coefficient component
 153 functions, whereupon we obtain an approximated AM and VCM by substitut-
 154 ing the resultant estimators into VCAM (2.2). In this way, we can estimate the
 155 varying-coefficient component functions and additive component functions.

156 **Step I:** Initial M-estimators of varying-coefficient component functions

157 In this step, we assume that B-spline bases have \bar{h}_C and \bar{h}_A interior knots for
 158 α_k and β_k , respectively. Using the tensor product of B-spline bases, the bivariate
 159 function $g_k(t, x_k) = \alpha_k(t)\beta_k(x_k)$ can be approximated as $g_k(t, x_k) \approx \gamma_k^T \mathcal{T}_k(t, x_k)$,
 160 where γ_k is a $\{(q + \bar{h}_C)(q + \bar{h}_A - 1)\}$ -vector. Assume that $\hat{\gamma} = (\hat{\gamma}_0^T, \dots, \hat{\gamma}_p^T)^T$ is
 161 determined by the following minimization problem:

$$\min_{\gamma} \sum_{i=1}^n \frac{1}{n_i} \sum_{j=1}^{n_i} \rho \left(y_{ij} - \gamma_0^T \mathbf{b}_C(t_{ij}) - \sum_{k=1}^p \gamma_k^T \mathcal{T}_k(t_{ij}, x_{ijk}) \right), \quad (2.3)$$

162 where ρ is a given loss function and $\gamma = (\gamma_0^T, \dots, \gamma_p^T)^T$.

163 For each given k , finding a point (t_{k0}, x_{k0}) such that $g_k(t_{k0}, x_{k0}) \neq 0$, then

164 $\xi_k(t|t_{k0}) = \frac{g_k(t, x_{k0})}{g_k(t_{k0}, x_{k0})} = \frac{\alpha_k(t)}{\alpha_k(t_{k0})}$ is well-defined and depends on the selection of t_{k0} .

165 Denoting $\hat{g}_k(t, x_k) = \hat{\gamma}_k^T \mathcal{T}_k(t, x_k)$, we can approximate $\xi_k(t|t_{k0})$ as $\hat{\xi}_k(t|t_{k0}, x_{k0}) =$

166 $\frac{\hat{g}_k(t, x_{k0})}{\hat{g}_k(t_{k0}, x_{k0})}$, which depends on the selection of t_{k0} and x_{k0} . In combination with

167 the identification conditions of α_k , we obtain the spline-based initial M-estimator

168 of $\alpha_k (k = 0, \dots, p)$ as

$$\hat{\alpha}_{0,I}(t) = \hat{\gamma}_0^T \mathbf{b}_C(t), \quad \hat{\alpha}_{k,I}(t|t_{k0}, x_{k0}) = \frac{\hat{\xi}_k(t|t_{k0}, x_{k0})}{\int_0^1 \hat{\xi}_k(t|t_{k0}, x_{k0}) dt}, \quad (2.4)$$

169 where the subscript ‘‘I’’ denotes the initial estimator of α_k .

170 **Step II:** M-estimators of additive component functions

171 Substituting (2.4), the initial M-estimator of α_k obtained in the Step-I es-

172 timation, into VCAM (2.2), we obtain the approximated AM $y_{ij} \approx \hat{\alpha}_{0,I}(t_{ij}) +$

173 $\sum_{k=1}^p \hat{\alpha}_{k,I}(t_{ij}|t_{k0}, x_{k0})\beta_k(x_{ijk}) + \varepsilon_{ij}$, which gives a spline-based M-estimator of

174 β_k . Denote the number of interior knots of the B-spline basis as K_A . Let

175 $\boldsymbol{\theta} = (\theta_1^\tau, \dots, \theta_p^\tau)^\tau$, with θ_k being a $(q + K_A - 1)$ -vector, so that $\hat{\boldsymbol{\theta}} = (\hat{\theta}_1^\tau, \dots, \hat{\theta}_p^\tau)^\tau$

176 minimizes the following problem:

$$\sum_{i=1}^n \frac{1}{n_i} \sum_{j=1}^{n_i} \rho\left(y_{ij} - \hat{\alpha}_{0,I}(t_{ij}) - \sum_{k=1}^p \hat{\alpha}_{k,I}(t_{ij}|t_{k0}, x_{k0})\theta_k^\tau \mathbf{B}_{k,A}(x_{ijk})\right). \quad (2.5)$$

177 Then the spline-based M-estimators $\hat{\beta}_k$, $k = 1, \dots, p$ of the additive component

178 functions are given by

$$\hat{\beta}_k(x_k) = \check{\beta}_k(x_k) - \frac{1}{N} \sum_{i=1}^n \sum_{j=1}^{n_i} \check{\beta}_k(x_{ijk}), \quad (2.6)$$

179 where $\check{\beta}_k(x_k) = \hat{\theta}_k^\tau \mathbf{B}_{k,A}(x_k)$ and $N = \sum_{i=1}^n n_i$.

180 **Step III:** Updated M-estimators of varying-coefficient component functions

181 Inserting (2.6) into (2.2), we obtain an approximated VCM, namely $y_{ij} \approx$

182 $\alpha_0(t_{ij}) + \sum_{k=1}^p \alpha_k(t_{ij})\hat{\beta}_k(x_{ijk}) + \varepsilon_{ij}$. Let K_C be the number of interior knots of

183 the B-spline basis fitting α_k . Denote $\mathbf{h} = (h_0^\tau, \dots, h_p^\tau)^\tau$, with h_k being a $(q + K_C)$ -

184 vector, so that $\hat{\mathbf{h}} = (\hat{h}_0^\tau, \dots, \hat{h}_p^\tau)^\tau$ minimizes

$$\sum_{i=1}^n \frac{1}{n_i} \sum_{j=1}^{n_i} \rho\left(y_{ij} - h_0^\tau \mathbf{b}_C(t_{ij}) - \sum_{k=1}^p \hat{\beta}_k(x_{ijk})h_k^\tau \mathbf{b}_C(t_{ij})\right). \quad (2.7)$$

185 Then the updated M-estimators of α_k , $k = 0, \dots, p$ are given by

$$\hat{\alpha}_0(t) = \hat{h}_0^\tau \mathbf{b}_C(t), \quad \hat{\alpha}_k(t) = \frac{\hat{h}_k^\tau \mathbf{b}_C(t)}{\int_0^1 \hat{h}_k^\tau \mathbf{b}_C(t) dt}.$$

186 Among the common convex loss functions are the quadratic function $\rho(u) =$
187 u^2 , the check function $\rho(u) = |u| + (2\tau - 1)u$ with $\tau \in (0, 1)$, and the Huber
188 function $\rho(u) = 0.5u^2 \mathbf{I}_{|u| < \delta}$, where δ is a pre-specified threshold value and \mathbf{I}_A
189 denotes the indicator function of a nonempty set A . Our method also allows for
190 a non-convex loss function, such as those of Hampel and Tukey. Note that the
191 estimation method proposed herein has a wide range of applications because the
192 spline approximations in the three estimation steps are valid for both sparse data
193 and dense data, simultaneously allowing the covariates to depend on the obser-
194 vation time. A simulation example given in Section S1.3 of the supplementary
195 material also compares our estimation method with that of Zhang and Wang
196 (2015) when the covariates are independent of the observation time. Table 6 in
197 the supplementary material shows that our estimators are superior to Zhang's
198 estimators for sparse data and a small proportion of outliers, and both perform
199 similarly for dense data with a normal error distribution.

200 3. Asymptotic Results

201 In this section, we construct the consistency and asymptotic normality of
202 the proposed M-estimators. Note that the asymptotic properties are considered
203 for sparse data and dense data in a unified framework, which can be viewed

204 as a generalization of Zhang and Wang (2016) to a VCAM. The necessary
 205 assumptions for deriving the asymptotic results are given in the Appendix.

206 **3.1. Consistency of three-step M-estimators**

207 Let $\bar{N}_H = (\sum_{i=1}^n n_i^{-1}/n)^{-1}$ be the harmonic average of sequence $\{n_i\}$, and
 208 let $\bar{h} = \bar{h}_A \vee \bar{h}_C$ be the maximum of \bar{h}_A and \bar{h}_C . Denote $\mathcal{J} = \{(\mathbf{x}_{ij}, t_{ij}) :$
 209 $, i = 1, \dots, n; j = 1, \dots, n_i\}$. Theorem 1 presents the rate of convergence for the
 210 additive component function β_k in the sense of the L_2 norm and mean squared
 211 errors (MSEs).

212 **Theorem 1.** *Under Assumptions A1–A5, M1 and M2, or N1 and N2, if $\bar{h} =$*
 213 *$O(K_A)$, $\bar{h}^2 K_A^{2r} = o(n\bar{N}_H)$, $K_A^2 = o(n\bar{N}_A)$, $K_A^{2r}/n \rightarrow C_1$, $K_A^{2r+1}/(n\bar{N}_H) \rightarrow C_2$,*
 214 *and $K_A/\bar{N}_H \rightarrow C_3$, where $0 \leq C_1 < \infty$, $0 \leq C_2, C_3 \leq \infty$, then we have the*
 215 *convergence rates*

$$\left\| \hat{\beta}_k - \beta_k \right\|_{L_2}^2 = O_p \left(K_A^{-2r} + \frac{K_A}{n\bar{N}_H} + \frac{1}{n} \right)$$

216 *in the L_2 -norm sense and*

$$\frac{1}{N} \sum_{i=1}^n \sum_{j=1}^{n_i} [\hat{\beta}_k(x_{ijk}) - \beta_k(x_{ijk})]^2 = O_p \left(K_A^{-2r} + \frac{K_A}{n\bar{N}_H} + \frac{1}{n} \right)$$

217 *in the MSE sense.*

218 **Remark 2.** It is easy to show that

219 (i) $\frac{1}{n} = o\left(\frac{K_A}{n\bar{N}_H}\right)$ if $\bar{N}_H/n^{2r} \rightarrow 0$ and $K_A \asymp (n\bar{N}_H)^{\frac{1}{2r+1}}$;

220 (ii) $\frac{1}{n} \asymp \frac{K_A}{n\bar{N}_H}$ if $\bar{N}_H/n^{2r} \rightarrow C$ and $K_A \asymp n^{\frac{1}{2r}}$;

221 (iii) $\frac{K_A}{n\bar{N}_H} = o(\frac{1}{n})$ if $\bar{N}_H/n^{2r} \rightarrow \infty$ and $K_A = o(n^{\frac{1}{2r}})$.

222 That is, the order of the variance term $\frac{K_A}{n\bar{N}_H} + \frac{1}{n}$ has either a parametric rate
 223 of convergence $\frac{1}{n}$ or a nonparametric rate of convergence $\frac{K_A}{n\bar{N}_H}$ according to the
 224 magnitude of \bar{N}_H/n^{2r} .

225 Theorem 2 is the analogue of Theorem 1 for the varying-coefficient function
 226 α_k .

227 **Theorem 2.** Under Assumptions A1–A5, M1 and M2, or N1 and N2, if $K_A K_C^{2r} =$
 228 $o(n\bar{N}_H)$, $K_A = O(K_C)$ or $K_A = o(K_C)$, $K_C^{2r}/n \rightarrow C_1$, $K_C^{2r+1}/(n\bar{N}_H) \rightarrow C_2$, and
 229 $K_C/\bar{N}_H \rightarrow C_3$, where $0 \leq C_1 < \infty$, $0 \leq C_2, C_3 \leq \infty$, then we have

$$\|\hat{\alpha}_k - \alpha_k\|_{L_2}^2 = O_p\left(K_C^{-2r} + \frac{K_C}{n\bar{N}_H} + \frac{1}{n}\right)$$

230 in the L_2 -norm sense and

$$\frac{1}{N} \sum_{i=1}^n \sum_{j=1}^{n_i} [\hat{\alpha}_k(t_{ij}) - \alpha_k(t_{ij})]^2 = O_p\left(K_C^{-2r} + \frac{K_C}{n\bar{N}_H} + \frac{1}{n}\right)$$

231 in the MSE sense.

232 A remark similar to that made regarding Theorem 1 can be made for M-
 233 estimators of varying-coefficient functions. Based upon these statements, we say
 234 that the data are

- 235 • sparse if $\bar{N}_H/n^{\frac{1}{2r}} \rightarrow 0$, which yields a nonparametric rate; or
- 236 • dense if $\bar{N}_H/n^{\frac{1}{2r}} \rightarrow C$ with $0 < C \leq \infty$, which yields a parametric rate.

237 It can be said that we generalize the way of splitting sparse data and dense data
 238 in the sense that our conclusions reduce to those of Zhang and Wang (2016)
 239 when $r = 2$.

240 3.2. Asymptotic normality of three-step M-estimators

In this subsection, we present the asymptotic distribution of M-estimators.

First, we introduce the following notation:

$$\begin{aligned}
 W_{n,A} &= \sum_{i=1}^n \frac{1}{n_i} \sum_{j=1}^{n_i} \varpi(t_{ij}) \Psi_{ij} \Psi_{ij}^\tau, & U_{n,A} &= \sum_{i=1}^n \Psi_i^\tau G_i \Psi_i / n_i^2, \\
 W_{n,C} &= \sum_{i=1}^n \frac{1}{n_i} \sum_{j=1}^{n_i} \varpi(t_{ij}) \Phi_{ij} \Phi_{ij}^\tau, & U_{n,C} &= \sum_{i=1}^n \Phi_i^\tau G_i \Phi_i / n_i^2,
 \end{aligned}$$

where

$$\begin{aligned}
 \Psi_i &= \{\Psi_{i1}, \dots, \Psi_{in_i}\}^\tau, \quad \Psi_{ij} = \{\psi_1(x_{ij1})^\tau, \dots, \psi_p(x_{ijp})^\tau\}^\tau, \quad \psi_k(x_{ijk}) = \alpha_k(t_{ij}) \mathbf{B}_{k,A}(x_{ijk}), \\
 \Phi_i &= \{\Phi_{i1}, \dots, \Phi_{in_i}\}^\tau, \quad \Phi_{ij} = \{1, \beta_1(x_{ij1}), \dots, \beta_p(x_{ijp})\}^\tau \otimes \mathbf{b}_C(t_{ij}).
 \end{aligned}$$

241 Theorem 3 presents the asymptotic distribution for the additive function β_k .

242 **Theorem 3.** Under the conditions of Theorem 1, if $K_A^{2r} \tilde{K}_A/n \rightarrow \infty$,

$$\frac{\max \left(K_A^{3/2} \sum_{i=1}^n 1/n_i^2, K_A^{1/2} \sum_{i=1}^n (n_i - 1)/n_i^2, \sum_{i=1}^n (n_i - 1)(n_i - 2)/n_i^2 \right)}{\left(\sum_{i=1}^n \frac{1}{n_i} (K_A - 1) + n \right)^{3/2}} \rightarrow 0$$

243 and the largest eigenvalue of $K_A \mathbf{B}_{k,A}(x) \mathbf{B}_{k,A}(x)^\tau$ is bounded, then given \mathcal{J} , it
 244 follows that $\hat{\beta}_k(x) - \beta_k(x) \xrightarrow{D} N(0, D_{n,A}(x))$, where

$$D_{n,A}(x) = A_k(x)^\tau W_{n,A}^{-1} U_{n,A} W_{n,A}^{-1} A_k(x), \quad (3.1)$$

245 and $A_k(x) = \{\mathbf{0}, \dots, \mathbf{B}_{k,A}^\tau(x), \dots, \mathbf{0}\}^\tau$ is a $\{pJ_A\}$ -dimensional vector with $\mathbf{B}_{k,A}(x)$
 246 at its $\{(k-1)J_A\}$ th to $\{kJ_A\}$ th positions and zero at the rest.

247 Theorem 4 is the analogue of Theorem 3 for the varying-coefficient function
 248 α_k .

249 **Theorem 4.** Under the conditions of Theorem 2, if $K_C^{2r} \tilde{K}_C/n \rightarrow \infty$,

$$\frac{\max\left(K_C^{3/2} \sum_{i=1}^n 1/n_i^2, K_C^{1/2} \sum_{i=1}^n (n_i-1)/n_i^2, \sum_{i=1}^n (n_i-1)(n_i-2)/n_i^2\right)}{\left(\sum_{i=1}^n \frac{1}{n_i} (K_C-1) + n\right)^{3/2}} \rightarrow 0 \quad (3.2)$$

250 and the largest eigenvalue of $K_C \mathbf{b}_C(t) \mathbf{b}_C(t)^\tau$ is bounded, then given \mathcal{J} , it follows
 251 that $\hat{\alpha}_k(t) - \alpha_k(t) \xrightarrow{D} N(0, D_{n,C}(t))$, where

$$D_{n,C}(t) = C_k(t)^\tau W_{n,C}^{-1} U_{n,C} W_{n,C}^{-1} C_k(t), \quad (3.3)$$

252 and $C_k(t) = \{\mathbf{0}, \dots, \mathbf{b}_C^\tau(t), \dots, \mathbf{0}\}^\tau$ is a $\{(p+1)J_C\}$ -dimensional vector with $\mathbf{b}_C(t)$
 253 at its $\{kJ_C\}$ th to $\{(k+1)J_C\}$ th positions and zero at the rest.

Now, we build a consistent estimate of asymptotic variance given in (3.1)
 and (3.3). Let $\hat{G}_i = \phi(\hat{\varepsilon}_i) \phi(\hat{\varepsilon}_i)^\tau$ with $\phi(\hat{\varepsilon}_i) = \{\phi(\hat{\varepsilon}_{i1}), \dots, \phi(\hat{\varepsilon}_{in_i})\}^\tau$ and $\hat{\varepsilon}_{ij} =$

$y_{ij} - \hat{\alpha}_0(t_{ij}) - \sum_{k=1}^p \hat{\alpha}_k(t_{ij}) \hat{\beta}_k(x_{ijk})$. Set

$$\begin{aligned} \hat{W}_{n,A} &= \sum_{i=1}^n \frac{1}{n_i} \sum_{j=1}^{n_i} \varpi(t_{ij}) \hat{\Psi}_{ij} \hat{\Psi}_{ij}^\tau, & \hat{U}_{n,A} &= \sum_{i=1}^n \hat{\Psi}_i \hat{G}_i \hat{\Psi}_i / n_i^2, \\ \hat{W}_{n,C} &= \sum_{i=1}^n \frac{1}{n_i} \sum_{j=1}^{n_i} \varpi(t_{ij}) \hat{\Phi}_{ij} \hat{\Phi}_{ij}^\tau, & \hat{U}_{n,C} &= \sum_{i=1}^n \hat{\Phi}_i \hat{G}_i \hat{\Phi}_i / n_i^2, \end{aligned}$$

where $\hat{\Psi}_{ij}$ and $\hat{\Phi}_{ij}$ are the counterparts of Ψ_{ij} and Φ_{ij} , respectively, after replacing α_k with $\hat{\alpha}_{k,I}$ and replacing β_k with $\hat{\beta}_k$. Then, the natural estimates of $D_{n,A}(x)$ and $D_{n,C}(t)$ are

$$\begin{aligned} \hat{D}_{n,A}(x) &= A_k(x)^\tau \hat{W}_{n,A}^{-1} \hat{U}_{n,A} \hat{W}_{n,A}^{-1} A_k(x) \quad \text{and} \\ \hat{D}_{n,C}(t) &= C_k(t)^\tau \hat{W}_{n,C}^{-1} \hat{U}_{n,C} \hat{W}_{n,C}^{-1} C_k(t). \end{aligned}$$

254 Theorem 5 shows that the estimates of asymptotic variances are consistent.

255 **Theorem 5.** Suppose that $\sup_{t \in [0,1]} \mathbb{E}(\phi^4(\varepsilon_{ij}) | t_{ij} = t) < \infty$.

256 (i) Under the conditions of Theorem 3, if $K_A = o(\bar{h}^r)$, $K_A^2 = o(n\bar{N}_H)$, $\max_i n_i K_A^2$
 257 $= o(n\bar{N}_H)$, and $K_A^4 \max_i n_i = o(n^4)$, then it holds that $\hat{D}_{n,A}(x) \xrightarrow{P} D_{n,A}(x)$.

258 (ii) Under the conditions of Theorem 4, if $K_C = o(K_A)$, $K_C^2 = o(n\bar{N}_H)$,
 259 $K_C^2 \max_i n_i = o(n\bar{N}_H)$, and $K_C^4 \max_i n_i = o(n^4)$, then it holds that $\hat{D}_{n,C}(x) \xrightarrow{P}$
 260 $D_{n,C}(x)$.

261 In combination with Theorems 3–5, the $(1 - \alpha)\%$ confidence intervals of
 262 univariate component functions are given by

$$\hat{\alpha}_k(t) \pm z_{\alpha/2} \{\hat{D}_{n,C}(t)\}^{1/2} \quad \text{and} \quad \hat{\beta}_k(x) \pm z_{\alpha/2} \{\hat{D}_{n,A}(x)\}^{1/2}. \quad (3.4)$$

263 **3.3. Quantile regression**

264 Let $0 < \tau < 1$ and loss function $\rho(u) = |u| + (2\tau - 1)u$, then the proposed
265 M-estimators reduce to τ -th quantile estimates. Denote $\hat{\alpha}_{k,\tau}(t)$ and $\hat{\beta}_{k,\tau}(x)$ as
266 the τ -th quantile estimates of α_k and β_k , respectively. We impose the following
267 additional assumptions.

268 (Q1) $P(\varepsilon_{ij} \leq 0 | \mathbf{x}_{ij}, t_{ij}) = \tau$.

269 (Q2) There exist positive constants c_5 and C_6 such that the conditional density
270 function $g(x|t)$ of ε_{ij} given $t_{ij} = t$ satisfies $|g(x|t) - g(0|t)| \leq C_6|x|$ for
271 all $x \in [-c_5, c_5]$ and $t \in [0, 1]$, and $g(0|t)$ is bounded away from zero and
272 infinity uniformly over $[0, 1]$.

273 Noting that $\rho(u)$ is convex and $\phi(u) = \rho'(u) = 2\tau I(u > 0) + 2(\tau - 1)I(u <$
274 $0)$, it is easy to show that Assumption M2 holds. If Assumption Q1 holds,
275 then $E\phi(\varepsilon_{ij}) = 0$ and Assumption M1 holds with $\varpi(t) = 2g(0|t)$. Employing
276 Theorems 1 and 2, we obtain the following corollary.

277 **Corollary 1.** *Suppose that conditions Q1 and Q2 hold.*

278 • *Under the conditions of Theorem 1, we have*

$$\left\| \hat{\beta}_{k,\tau} - \beta_k \right\|_{L_2}^2 = O_p \left(K_A^{-2r} + \frac{K_A}{n\bar{N}_H} + \frac{1}{n} \right).$$

279 • Under the conditions of Theorem 2, we have

$$\|\hat{\alpha}_{k,\tau} - \alpha_k\|_{L_2}^2 = O_p\left(K_C^{-2r} + \frac{K_C}{n\bar{N}_H} + \frac{1}{n}\right).$$

280 **Remark 3.** Let $\varpi(t) = 2g(0|t)$ in $W_{n,A}$ and $W_{n,C}$. If conditions Q1 and Q2 hold,
 281 then we can present the asymptotic distributions of $\hat{\beta}_{k,\tau}(x)$ and $\hat{\alpha}_{k,\tau}(t)$ under the
 282 conditions of Theorems 3 and 4, respectively.

283 4. Model Identification Procedure

284 The VCAM (2.1) is a flexible nonparametric regression method. However,
 285 parsimony is always preferable when several potential options are available. To
 286 this end, we propose a model identification strategy based on the penalized M-
 287 estimators to identify additive terms and varying-coefficient terms.

288 The assumption of continuous covariates means that $X_k \neq 0$ almost surely
 289 for $k = 1, \dots, p$, and model (2.2) can be rewritten as

$$y_{ij} = \alpha_0(t_{ij}) + \sum_{k=1}^p x_{ijk} \alpha_k(t_{ij}) \beta_k^*(x_{ijk}) + \varepsilon_{ij},$$

where $\beta_k^*(x) = \beta_k(x)/x$. Employing the tensor product of B-spline bases, the
 bivariate function $g_k^*(t, x_k) = \alpha_k(t) \beta_k^*(x_k)$ can be approximated as

$$\begin{aligned} g_k^*(t, x_k) &\approx \{1, \mathbf{B}_{k,AP}^\tau(x_k)\} \otimes \{1, \mathbf{B}_{CP}^\tau(t)\} \eta_k \\ &= \eta_{00,k} + \eta_{0,0,k}^\tau \mathbf{B}_{CP}(t) + \eta_{01,k} B_{k1}(x_k) + \eta_{1,1,k}^\tau B_{k1}(x_k) \otimes \mathbf{B}_{CP}(t) \\ &\quad + \cdots + \eta_{0J_{AP},k} B_{kJ_{AP}}(x_k) + \eta_{J_{AP},k}^\tau B_{kJ_{AP}}(x_k) \otimes \mathbf{B}_{CP}(t), \end{aligned}$$

290 where $\eta_k = \{\eta_{00,k}, \eta_{0,k}^\tau, \eta_{01,k}, \eta_{1,k}^\tau, \dots, \eta_{0J_{AP},k}, \eta_{J_{AP},k}^\tau\}^\tau$, $\eta_{j,k} = \{\eta_{1j,k}, \dots, \eta_{J_{CP}j,k}\}^\tau$,
 291 and J_{AP} and J_{CP} are the cardinalities of the B-spline bases $\mathbf{B}_{k,AP}(x_k)$ and $\mathbf{B}_{CP}(t)$
 292 for β_k and α_k , respectively, in the model identification procedure.

293 Let $M_k(t, x_k) = \{0, \mathbf{B}_{CP}^\tau(t), 0, B_{k1}(x_k) \otimes \mathbf{B}_{CP}^\tau(t), \dots, 0, B_{kJ_{AP}}(x_k) \otimes \mathbf{B}_{CP}^\tau(t)\}^\tau$
 294 and $F_k(t, x_k) = \{0_{J_{CP}+1}^\tau, \mathbf{B}_{k,AP}^\tau(x_k) \otimes (1, \mathbf{B}_{CP}^\tau(t))\}^\tau$, where 0_l denotes the l -vector
 295 of zeros. Then, we can say that g_k reduces to

- 296 • an additive term if and only if $\eta_k^\tau M_k(t, x_k) = 0$ and
- 297 • a varying-coefficient term if and only if $\eta_k^\tau F_k(t, x_k) = 0$

298 for any $(t, x) \in [0, 1] \times [a_k, b_k]$, where $[a_k, b_k]$ is the domain of $\beta_k(\cdot)$.

299 We now propose a regularized M-estimation method by penalizing the L_2
 300 norm of $M_k^\tau \eta_k$ and $F_k^\tau \eta_k$ for $k = 1, \dots, p$. Denote the numbers of interior knots
 301 for α_k and β_k in the model identification procedure as \tilde{h}_{CP} and \tilde{h}_{AP} , respectively.
 302 Let $\boldsymbol{\eta} = (\eta_0^\tau, \dots, \eta_p^\tau)^\tau$, with η_0 being a $\{q + \tilde{h}_{CP}\}$ -vector and $\eta_k (k = 1, \dots, p)$ being
 303 a $\{(q + \tilde{h}_{CP})(q + \tilde{h}_{AP} - 1)\}$ -vector. Suppose that $\hat{\boldsymbol{\eta}} = (\hat{\eta}_0^\tau, \dots, \hat{\eta}_p^\tau)^\tau$ minimizes the
 304 following problem:

$$\begin{aligned} & \sum_{i=1}^n \frac{1}{n_i} \sum_{j=1}^{n_i} \rho \left(y_{ij} - \eta_0^\tau \mathbf{b}_C(t_{ij}) - \sum_{k=1}^p x_{ijk} \{1, \mathbf{B}_{k,AP}^\tau(x_k)\} \otimes \{1, \mathbf{B}_{CP}^\tau(t)\} \eta_k \right) \\ & + n \sum_{k=1}^p p_{\lambda_1} (\|M_k^\tau \eta_k\|_{L_2}) + n \sum_{k=1}^p p_{\lambda_2} (\|F_k^\tau \eta_k\|_{L_2}). \end{aligned} \quad (4.1)$$

305 The product term $\alpha_k(t)\beta_k(x_k)$ in (1.1) then becomes an additive term if $\|M_k^\tau \hat{\eta}_k\|_{L_2}$
 306 is close to zero (e.g., no larger than 10^{-4}) and becomes a varying-coefficient term

307 if $\|F_k^T \hat{\eta}_k\|_{L_2}$ is close to zero.

308 There are various ways to specify the penalty function $p_\lambda(\cdot)$ (Tibshirani ,
309 1996; Fan and Li , 2001; Zou , 2006). We adopt the smoothly clipped absolute
310 deviation (SCAD) penalty function and use the locally quadratic approximation
311 (LQA) algorithm proposed by Fan and Li (2001).

312 Let \mathcal{I}_A and \mathcal{I}_V be the index sets of additive terms and varying-coefficient
313 terms in VCAM (2.1), respectively. Denote $\varrho_n = \bar{h}_P^{-r} + \sqrt{\kappa_P/n}$ with $\bar{h}_P =$
314 $\bar{h}_{AP} \wedge \bar{h}_{CP}$ and $\kappa_P = \bar{h}_P^2 / \bar{N}_H$.

315 Theorem 6 demonstrates the consistency of the model identification proce-
316 dure.

317 **Theorem 6.** *Suppose that Assumptions A1–A5, M1 and M2, or N1 and N2*
318 *hold.*

319 (i) *If $\lambda_1 \rightarrow 0$, $\sqrt{\varrho_n}/\lambda_1 \rightarrow 0$, and $\liminf_{n \rightarrow \infty} \liminf_{w \rightarrow 0^+} p'_{\lambda_1}(w)/\lambda_1 = 1$, then*

320 $M_k^T(t, x_k) \hat{\eta}_k = 0 \forall k \in \mathcal{I}_A$ *with probability approaching unity.*

321 (ii) *If $\lambda_2 \rightarrow 0$, $\sqrt{\varrho_n}/\lambda_2 \rightarrow 0$, and $\liminf_{n \rightarrow \infty} \liminf_{w \rightarrow 0^+} p'_{\lambda_2}(w)/\lambda_2 = 1$, then*

322 $F_k^T(t, x_k) \hat{\eta}_k = 0 \forall k \in \mathcal{I}_V$ *with probability approaching unity.*

323 5. Implementation Issues

324 In this section, we address several practical problems regarding the selection
325 of smoothing parameters and tuning parameters in our methods. As is common

326 practice in the spline literature, we select the number of interior knots via a
327 data-driven method [namely the Bayes information criterion (BIC)] and position
328 them at equal intervals on the sample quantiles.

- 329 • Selecting the optimal number of interior knots (\hat{h}_C, \hat{h}_A) .

330 The optimal number of interior knots (\hat{h}_C, \hat{h}_A) in the Step-I estimation
331 minimizes the following BIC function:

$$\text{BIC}_1(\hat{h}_C, \hat{h}_A) = \log \left(\frac{1}{n} \sum_{i=1}^n \frac{1}{n_i} \sum_{j=1}^{n_i} \rho(\hat{\sigma}_{ij,1}) \right) + \frac{\log N}{2N} N_1,$$

332 where $\hat{\sigma}_{ij,1} = y_{ij} - \hat{\gamma}_0^\top \mathbf{b}_C(t_{ij}) - \sum_{k=1}^p \hat{\gamma}_k^\top \mathcal{T}_k(t_{ij}, x_{ijk})$ and $N_1 = (q + \hat{h}_C)(1 +$
333 $p(q + \hat{h}_A - 1))$.

- 334 • Selecting the optimal number of interior knots (K_A, K_C) .

335 The optimal number of interior knots (\hat{K}_A, \hat{K}_C) in Steps II and III mini-
336 mizes

$$\text{BIC}_2(K_A, K_C) = \log \left(\frac{1}{n} \sum_{i=1}^n \frac{1}{n_i} \sum_{j=1}^{n_i} \rho(\hat{\sigma}_{ij,2}) \right) + \frac{\log N}{2N} N_2,$$

337 where $\hat{\sigma}_{ij,2} = y_{ij} - \hat{\alpha}_0(t_{ij}) - \sum_{k=1}^p \hat{\alpha}_k(t_{ij}) \hat{\beta}_k(x_{ijk})$ and $N_2 = p(q + K_A -$
338 $1) + (p + 1)(q + K_C)$.

- 339 • Selecting the optimal tuning parameters (λ_1, λ_2) .

340 We use the optimal number of interior knots (\hat{h}_C, \hat{h}_A) and the optimal
 341 tuning parameters $(\hat{\lambda}_1, \hat{\lambda}_2)$ that minimize the following BIC:

$$\text{BIC}_3(\lambda_1, \lambda_2) = \log \left(\frac{1}{n} \sum_{i=1}^n \frac{1}{n_i} \sum_{j=1}^{n_i} \rho(\hat{\sigma}_{ij,3}) \right) + \frac{\log N}{2N} N_3,$$

342 where $\hat{\sigma}_{ij,3} = y_{ij} - \hat{\eta}_0^\tau \mathbf{b}_{\text{CP}}(t_{ij}) - \sum_{k=1}^p x_{ijk} \{1, \mathbf{B}_{k,\text{AP}}^\tau(x_{ijk})\} \otimes \{1, \mathbf{B}_{\text{CP}}^\tau(t_{ij})\} \hat{\eta}_k$
 343 and $N_3 = m_L + \{q + \hat{h}_C\} \{m_C + 1\} + m_A \{q + \hat{h}_A - 1\} + \{q + \hat{h}_C\} \{q + \hat{h}_A -$
 344 $1\} \{p - m_L - m_C - m_A\}$, with m_L linear terms, m_A additive terms, and m_C
 345 varying-coefficient terms.

346 6. Numerical Studies

347 Simulation examples are used to investigate the finite-sample performance
 348 of the proposed three-step M-estimation method and model identification proce-
 349 dure. Empirical examples are then presented to illustrate the usefulness of our
 350 method in practice.

351 6.1. Simulation studies

352 **Example 1.** A VCAM with repeated measurements is generated as follows:

$$y_{ij} = \alpha_0(t_{ij}) + \alpha_1(t_{ij})\beta_1(x_{ij}) + w_i(t_{ij}) + e_{ij}, \quad i = 1, \dots, n; \quad j = 1, \dots, m,$$

353 where t_{ij} are i.i.d. copies from $U(0, 1)$, $x_{ij} = 0.8t_{ij}^2 + \eta_{ij}$ with η_{ij} drawn indepen-
 354 dently from $N(0, (1 + t_{ij})/(2 + t_{ij}))$. The subject-specific random trajectories
 355 $w_i(i = 1, \dots, n)$ are independent copies of a zero-mean stationary Gaussian pro-
 356 cess with covariance function $\gamma(u) = 0.35\theta^{|u|}$, where $\theta = 0$ and 0.5. The random

357 noise e_{ij} are i.i.d. from four error distributions, namely the normal distribution
358 $N(0, 0.2)$, the mixed normal distribution $0.95N(0, 0.2) + 0.05N(0, 12.5^2)$, and the
359 scaled t distributions of $0.5 \times t(2)$ and $0.2 \times t(1)$. The univariate component func-
360 tions are given by $\alpha_0(t) = \cos(2\pi t)$, $\alpha_1(t) = \{2t \sin(2\pi t) + 1\} / \int_0^1 \{2t \sin(2\pi t) +$
361 $1\} dt$, and $\beta_1(x) = 1.5 \sin(\pi x/2) - x(1-x) - E[1.5 \sin(\pi X/2) - X(1-X)]$.

362 Three loss functions are considered, namely the quadratic function $\rho_1(x) =$
363 x^2 , the absolute value function $\rho_2(x) = |x|$, and the Huber function $\rho_3(x) =$
364 $0.5x^2 \mathbf{I}_{|x| < \delta}$ with $\delta = 1.345$. We appraise the performance of the three-step M-
365 estimator through the MSE, which is defined as

$$\text{MSE}(g) = \frac{1}{nm} \sum_{i=1}^n \sum_{j=1}^m [\hat{g}(t_{ij}) - g(t_{ij})]^2,$$

where g is either α_k or β_k . To obtain an intuitive impression of the robustness
of the M-estimators, we define the weighted average squared error (WASE) as

$$\text{WASE} = \frac{1}{nm} \sum_{i=1}^n \sum_{j=1}^m \left\{ \frac{[\hat{\alpha}_0(t_{ij}) - \alpha_0(t_{ij})]^2}{[\text{range}(\alpha_0)]^2} + \frac{[\hat{\alpha}_1(t_{ij}) - \alpha_1(t_{ij})]^2}{[\text{range}(\alpha_1)]^2} + \frac{[\hat{\beta}_1(x_{ij}) - \beta_1(x_{ij})]^2}{[\text{range}(\beta_1)]^2} \right\},$$

366 where $\text{range}(f)$ denotes the range of a given function f .

367 For $n = 30$ and $m = 20$, based upon 500 Monte Carlo replications, Figure 1
368 shows the average WASE of three-step M-estimators with four error distribu-
369 tions and two types of intra-subject covariance structure. In this figure, 1, 2,
370 and 3 denote the least-squares estimator, the median estimator, and the Huber
371 estimator, respectively. We also compare the average MSE (AMSE) in Table 1
372 of the supplementary material. From the obtained results, we conclude that the

373 Huber estimator and the median estimator have similar performance whichever
374 error distribution is adopted. They are comparable to least-squares estimators
375 under normal error distributions and are superior to least-squares estimators un-
376 der non-normal error distributions. In addition, the influence of the intra-subject
377 covariance structure is insignificant.

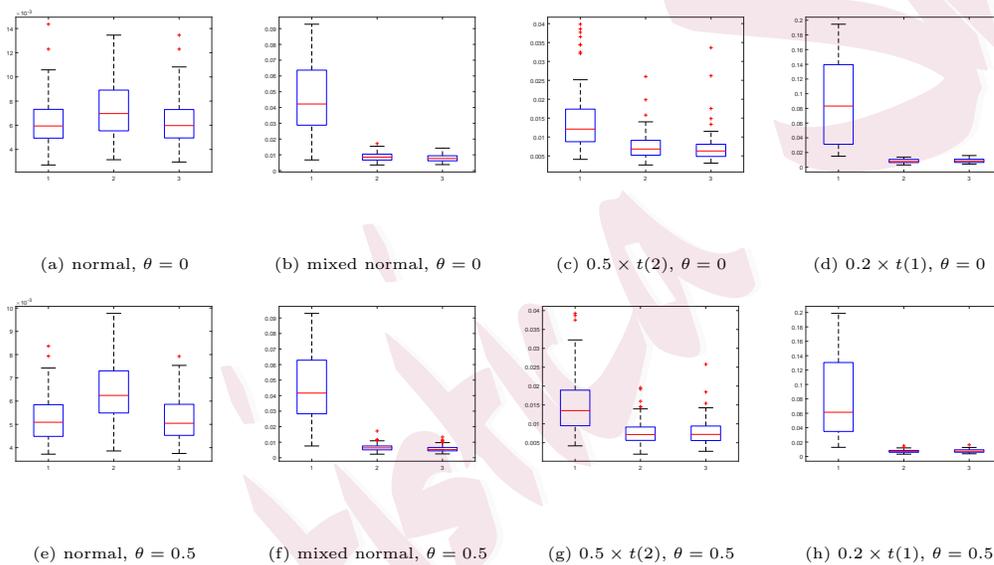


Figure 1: Boxplot for the average of WASEs (AWASE) based on 500 Monte Carlo replications. 1, 2, and 3 denote least-squares estimator, median estimator, and Huber estimator, respectively.

378 Furthermore, we investigate the graphical representation of the iterative Hu-
379 ber estimator under a mixed normal error distribution. Figure 2 shows the point-
380 wise 95% confidence intervals (CIs) of the Huber estimator based on the central

381 limit theorem (CLT) (dotted lines) and the 95% CIs based upon 500 wild boot-
382 strap samplings (dash-dotted lines). The true component function (solid line) and
383 Huber M-estimator (dashed line) are also given. The resulting figures indicate
384 that the difference between the two types of CIs is insignificant, which motivates
385 us to make the bold claim that the bootstrap method is sound. However, we do
386 not investigate the theoretical justification for that claim to avoid straying from
387 the aim of the present paper. Meanwhile, note that the true curves and the Hu-
388 ber estimators are very close and both fall into the 95% CIs, thereby indicating
389 the rationality of the proposed estimation method. Under a normal error distri-
390 bution, the least-squares-based CIs are shown in Figure 1 of the supplementary
391 material.

392 We also investigate the average experience coverage probability (AECp) of
393 the three-step M-estimator with a normal error distribution and a mixed normal
394 error distribution in Figures 2 and 3 of the supplementary material, respectively,
395 which show that the pointwise CLT-based CI performs well even in the presence
396 of a small proportion of outliers. In addition, Figure 4 in the supplementary
397 material also compares the AECp of the component functions under a more
398 general sampling plan, namely sparse observations for some subjects and dense
399 observations for other subjects. The results show that the more general sampling
400 plan and a small proportion of outliers have no significant influence on the AECp

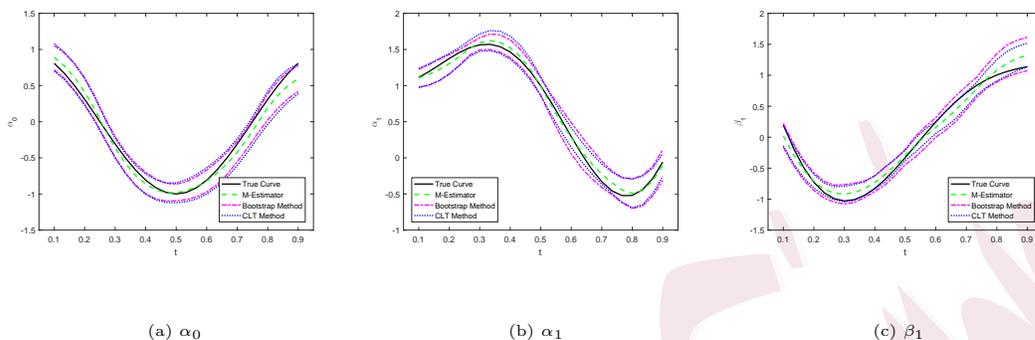


Figure 2: Three-step M-estimators under mixed normal error distribution.

Solid line: true component function; dashed line: three-step M-estimator; dotted lines: 95% CIs based on (3.4); dash-dotted lines: 95% CIs based on 500 wild bootstrap resamplings.

401 of the component functions.

402 Tables 2 and 3 in the supplementary material also compare the average
403 of MSE (AMSEs) of the iterative M-estimator under different combinations of
404 (n, m) with $n = 20, 40$ and $m = 20, 30$. We conclude that, as the total number
405 of observations grows, the AMSE decreases for a normal error distribution re-
406 gardless of which loss function is used. For non-normal error distributions, the
407 AMSEs of estimators based on loss functions ρ_2 and ρ_3 decrease, but the least-
408 squares estimator shows no significant improvement as the total observation size
409 grows.

410 Moreover, the numerical example considered in Section S1.2 of the supple-

mentary material investigates the finite-sample performance of the model identification procedure. As expected, the results given in Tables 4 and 5 of the supplementary material verify our asymptotic theories and demonstrate the robustness of the model identification.

6.2. Analysis of real data

Example 2. We now apply our method to the CD4 data from the Multicenter AIDS Cohort Study, which contain 2,376 observations from 369 men infected with HIV. Zhang, Park and Wang (2013) analyzed this dataset by means of the time-varying additive model $y_{ij} = \mu_0(t_{ij}) + \sum_{k=1}^2 \mu_k(t_{ij}, x_{ijk}) + w_{ij} + e_{ij}$. Following them, we also choose two covariates: X_1 (age) is the age at seroconversion (time-invariant variable) and X_2 (cesd) is the level of depression, which is recorded over time (in years).

Employing the separability testing proposed by Hu, Huang and You (2018), we obtain a p-value of 0.84, which means that the VCAM (2.2), a submodel of the time-varying additive model introduced by Zhang, Park and Wang (2013), is sufficient for this dataset. Under loss function ρ_3 in Example 1, we select optimal knots $(\hat{h}_C, \hat{h}_A, \hat{K}_C, \hat{K}_A) = (2, 2, 4, 3)$ according to the BIC given in Section 5 and obtain the optimal tuning parameters $(\hat{\lambda}_1, \hat{\lambda}_2) = (3.06, 1.56)$, which are selected from $[0.01, 5]$ with spacing 0.05. Based on the resulting optimal parameters, we obtain the penalized estimators and conclude that α_1 and α_2 are time-variant

431 and that β_1 and β_2 are nonlinear.

432 The Huber estimator and the 95% CIs of the univariate component functions
433 are presented in Figure 3, from which we conclude that the overall mean functions
434 α_0 and α_1 for X_1 (age) are monotonically decreasing and that α_2 for X_2 (cesd)
435 is a bimodal function. For a fixed time, the effect of age on the CD4 count
436 increases until around age = 12, after which it decreases. However, the effect of
437 depression on the CD4 count decreases rapidly before cesd = 5, then increases
438 until around cesd = 25, after which it decreases. The residuals plot in Figure 3(f)
439 shows that our regression method is appropriate for this dataset. Figure 6 in the
440 supplementary material also gives the estimated surfaces of the bivariate function
441 $g_k(t, x_k) = \alpha_k(t)\beta_k(x_k), k = 1, 2$.

442 **Example 3.** In this example, we consider a real diffusion-weighted imaging
443 dataset with $n = 213$ subjects collected from the NIH Alzheimer's Disease Neu-
444 roimaging Initiative (ADNI) study. The observed response process is a fractional
445 anisotropy (FA) curve at all 83 grid points along the skeleton of the midsagittal
446 corpus callosum. Here, we want to explore the relationship between FA (Y) and
447 three covariates, namely (i) the age of the subject (X_1), (ii) their educational
448 level (X_2), and (iii) the result of the ADNI Mini-Mental State Exam (X_3). Luo,
449 Zhu and Zhu (2016) and Li, Huang and Zhu (2017) analyzed this dataset us-
450 ing a single-index varying-coefficient model and a functional varying-coefficient

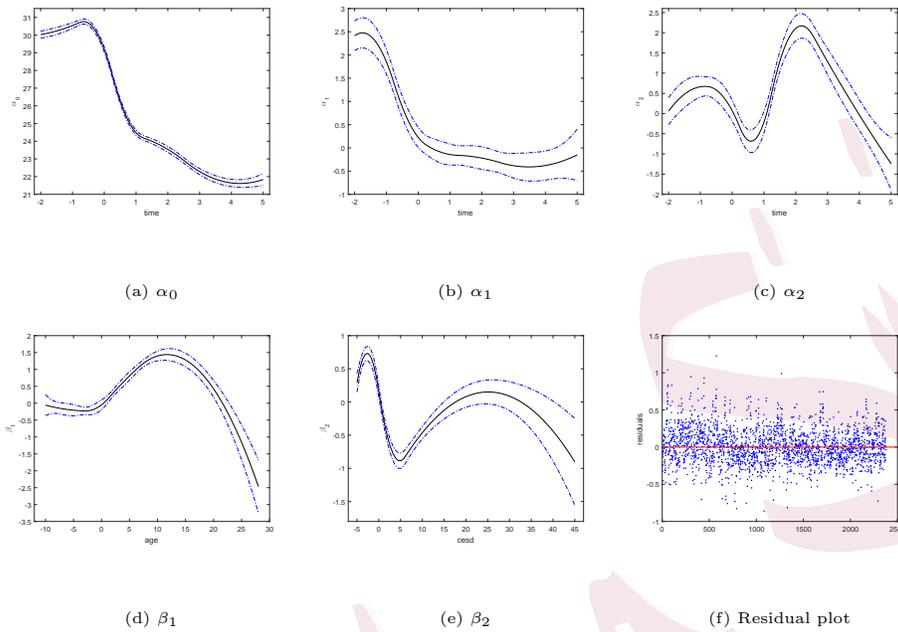


Figure 3: Three-step M-estimators for CD4 dataset. Solid line: three-step M-estimators; dash-dotted lines: 95% CIs based on (3.4); (f) plots scaled residuals relative to fitted values.

451 single-index model, respectively. The two models both assume linear covariate
452 effects with varying coefficients and/or nonlinear covariate effects only through
453 the linear combination of covariates with varying coefficients. However, the linear
454 effect is a somewhat strict constraint in practical applications, and we are inter-
455 ested in the function effect of each predictor on the response process, including
456 the linear effect as its special case. Therefore, we apply a VCAM to this dataset.

457 Employing the proposed model identification procedure, we claim that the

458 varying-coefficient functions are all time-variant and that the additive functions
459 are all nonlinear. Figure 4 shows the Huber estimators of the univariate compo-
460 nent functions and the 95% pointwise CIs based on (3.4). Figure 4(e)–(g) show
461 how the covariates affect the response process: the effect of age increases initially,
462 then decreases before the average age, and then increases thereafter; the effect
463 of educational level increases gently before the average educational level, then
464 decreases, and finally increases; the effect of the ADNI Mini-Mental State Exam
465 decreases until nearly the average value and then increases. The estimated bi-
466 variate functions $g_k(t, x_k) = \alpha_k(t)\beta_k(x_k)$, $k = 1, 2, 3$ are presented in Figure 10 of
467 the supplementary material, which shows the dynamic effects of the covariates.
468 The Q-Q plot shows that our regression method is appropriate for this dataset.

469 Analysis of the cigarette data mentioned in Section 1 shows that a reduced
470 VCAM is preferable. Details are given in Section S1.6 of the supplementary
471 material.

472 7. Concluding Remarks

473 The VCAM proposed by Zhang and Wang (2015) is a flexible structural
474 nonparametric regression method that includes the classical varying-coefficient
475 model and additive model as special cases. In this paper, we have considered an
476 M-type robust regression method for this VCAM to enable the analysis of lon-
477 gitudinal data and functional data, which may include sparse or dense repeated

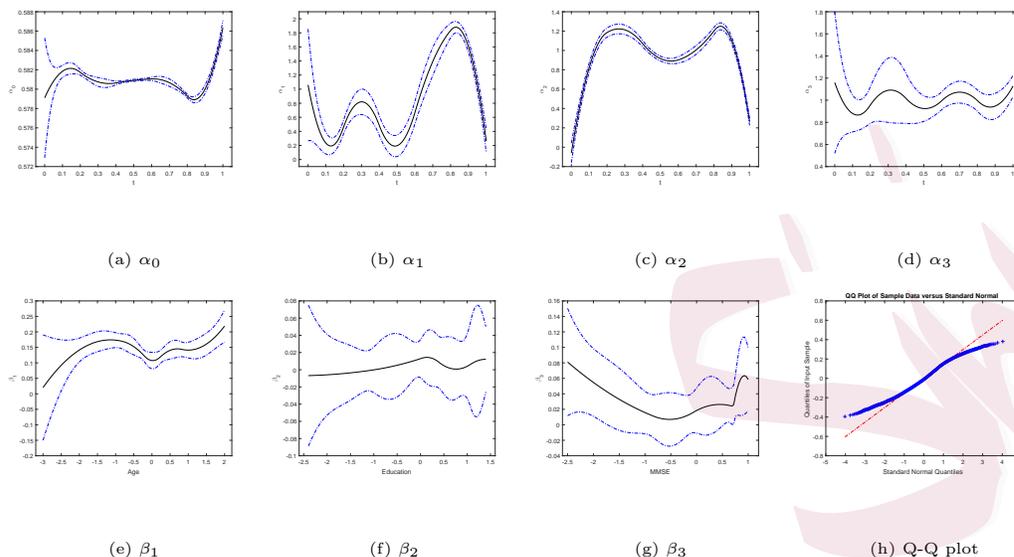


Figure 4: Estimated univariate component functions for ADNI data. Solid line: three-step M-estimator; dash-dotted lines: 95% CIs based on (3.4).

478 measurements for the selected subjects, and both the response and covariates
 479 may be smooth processes that are dependent on the observation time.

480 We have proposed spline-based three-step M-estimators for varying-coefficient
 481 component functions and additive component functions. The asymptotic prop-
 482 erties are considered for sparse and dense data in a unified framework, which
 483 yields a split of sparse data and dense data according to the relative order of n_i
 484 to n . Similar to Hu, Huang and You (2018), the proposed estimation method
 485 has the oracle property in that the iterative estimation procedure does not cause
 486 additional asymptotic errors.

487 To select as parsimonious a model as possible for the given real-life data, we
488 have also developed a model identification procedure based on the SCAD penalty
489 function. It has been shown that the proposed model identification method can
490 correctly select an additive term and a varying-coefficient term with probability
491 approaching unity.

492 Appendix

493 Let $C_r[a, b]$ be the space of all functions $m(x)$ defined on $[a, b]$ such that the
494 $(r - 1)$ -order derivative $m^{(r-1)}(\cdot)$ is continuous over $[a, b]$, and

$$|m^{(r-1)}(x) - m^{(r-1)}(x')| \leq C|x - x'|, \quad \forall x, x' \in [a, b],$$

495 where C is a positive constant. The necessary conditions for the asymptotic
496 results are listed below.

- 497 • Basic assumptions.

498 (A1) The time points $\{t_{ij}\}$ are independent copies of T , whose probabil-
499 ity density function $f_T(\cdot)$ is uniformly bounded away from zero and
500 infinity.

501 (A2) The marginal density function $f_k(\cdot)$ of X_k is uniformly bounded away
502 from zero and infinity over the support set S_k of X_k . The joint
503 density $f_{\mathbf{X}, T}(\mathbf{x}, t)$ of \mathbf{X} and T is uniformly bounded away from zero
504 and infinity on $(\mathbf{x}, t) \in \prod_{k=1}^p S_k \times [0, 1]$.

505 (A3) $\alpha_k \in C_r[0, 1]$ and $\beta_k \in C_r[a_k, b_k]$, where $1 \leq r \leq q$ and a_k, b_k are
506 finite real numbers for $k = 1, \dots, p$.

507 (A4) The function $\phi(\cdot) = \rho'(\cdot)$ satisfies $E[\phi(\varepsilon_{ij})|t_{ij} = t] = 0$ and $E[\phi^2(\varepsilon_{ij})|t_{ij} =$
508 $t] \leq C_1$ for any $t \in [0, 1]$, where C_1 is a positive constant.

509 (A5) There exists some positive constant $\tilde{\lambda}$ such that the smallest eigen-
510 value λ_{i1} of $G_i = E[\phi(\varepsilon_i)\phi(\varepsilon_i)^\tau|\mathcal{J}]$ satisfies $\lambda_{i1} \geq \tilde{\lambda} > 0$.

511 • Assumptions for convex loss function.

512 (M1) The loss function $\rho(\cdot)$ is convex, and there exist some function $\varpi(t)$
513 and positive constants c_1 and C_2 such that

$$|E[\phi(\varepsilon_{ij} + u)|t_{ij} = t] - \varpi(t)u| \leq C_2u^2$$

514 for any $|u| \leq c_1$ and $t \in [0, 1]$. Moreover, $\varpi(t)$ satisfies $0 < c_\varpi \leq$
515 $\min_{t \in [0, 1]} \varpi(t) \leq \max_{t \in [0, 1]} \varpi(t) \leq C_\varpi < \infty$.

516 (M2) There exist positive finite constants $c_2, C_3,$ and C_4 such that

$$E[\{\phi(\varepsilon_{ij} + u) - \phi(\varepsilon_{ij})\}^2|\mathcal{J}] \leq C_3|u|$$

517 and $|\phi(u + v) - \phi(v)| \leq C_4$ for any $|u| \leq c_2, t \in [0, 1]$, and $v \in \mathbb{R}$.

518 • Assumptions for non-convex loss function.

519 (N1) The function $\phi(\cdot)$ is continuous and has a derivative $\phi'(\cdot)$ almost
520 everywhere. Furthermore, $\phi_\varepsilon(t) = \mathbb{E}[\phi'(\varepsilon_{ij})|t_{ij} = t]$ is positive and
521 continuous at t .

522 (N2) $\mathbb{E}[\sup_{\|z\| \leq \delta} |\phi(\varepsilon_{ij} + z) - \phi(\varepsilon_{ij}) - \phi'(\varepsilon_{ij})z| | t_{ij} = t] = o(\delta)$ as $\delta \rightarrow 0$.

523 **Remark 4.** Assumptions A1 and A2 relate to the distributions of time points
524 t_{ij} and covariates \mathbf{x}_{ij} . Assumption A3 specifies the degree of smoothness of
525 varying-coefficient component functions and additive component functions. As-
526 sumptions A4, M1, and M2 are standard assumptions about the score function
527 ϕ of a convex loss function; see He, Zhu and Fung (2002); Tang and Cheng
528 (2008) for details. Assumptions N1 and N2 are necessary for a non-convex loss
529 function; see Fan and Jiang (2000); Jiang and Mack (2001).

530 Supplementary Material

531 The supplementary material includes additional numerical studies, an itera-
532 tive algorithm for penalized M-estimators, and proof techniques for deriving the
533 asymptotic results.

534 Acknowledgments

535 The authors are grateful to the Co-editor, an associate editor, and the refer-
536 ees for their constructive comments that have substantially improved the quality
537 of this article.

References

- 538
- 539 Bada, O. and Liebl D. (2012). phtt: Panel data analysis with heterogeneous time trends. R
540 package version 3 (2).
- 541 Carroll, R. J., Maity, A., Mammen, E., and Yu, K. (2009). Nonparametric additive regression
542 for repeatedly measured data. *Biometrika*, **96**, 383–398.
- 543 De Boor. (1978). *A Practical Guide to Splines*. Springer, New York.
- 544 Diggle P. J., Liang, K. Y., and Zeger, S. L. (1994). *Analysis of longitudinal data*. Oxford, U.K.:
545 Oxford University Press.
- 546 Fan, J. and Jiang, J. (2000). Variable bandwidth and one-step local M-estimator. *Science in*
547 *China*, **43**,65–81.
- 548 Fan, J., and Li, R. (2001), Variable selection via nonconcave penalized likelihood and its oracle
549 properties. *Journal of the American Statistical Association*, **96**, 1348–1360.
- 550 He, X., and Shi, P. (1994). Convergence rate of b-spline estimators of nonparametric conditional
551 quantile functions. *Journal of Nonparametric Statistics*, **3**, 299–308.
- 552 He, X., Zhu, Z., and Fung, W. (2002). Estimation in a semiparametric model for longitudinal
553 data with unspecified dependence structure. *Biometrika*, **89**, 579–590.
- 554 Hoover, D. R., Rice, J. A., Wu, C. O., and Yang, L.-P. (1998). Nonparametric smoothing
555 estimates of time-varying coefficient models with longitudinal data. *Biometrika*, **85**, 809–
556 822.

36REFERENCES

- 557 Hu, L. X., Huang, T., and You J. H. (2018). Estimation and Identification of a Varying-
558 Coefficient Additive Model for Locally Stationary Processes. *Journal of the American Sta-*
559 *tistical Association*, doi: 10.1080/01621459.2018.1482753.
- 560 Jiang, J. and Mack Y. (2001). Robust local polynomial regression for dependent data. *Statistica*
561 *Sinica*, **11**, 705–722.
- 562 Koenker, R., and Bassett, G. (1978). Regression quantiles. *Econometrica*. **46**, 33–50.
- 563 Li, Y. H., and Hsing, T. (2010). Uniform convergence rates for nonparametric regression and
564 principal component analysis in functional/longitudinal data. *The Annals of Statistics*, **38**,
565 3321 – 3351.
- 566 Li, J., Huang, C., Zhu, H., & for the Alzheimers Disease Neuroimaging Initiative (2017). A Func-
567 tional Varying-Coefficient Single-Index Model for Functional Response Data. *Journal of the*
568 *American Statistical Association*, **112**, 1169-1181, doi: 10.1080/01621459.2016.1195742.
- 569 Luo, X., Zhu, L., and Zhu, H. (2016). Single-index varying coefficient model for functional
570 responses. *Biometrics*, **72**, 1275–1284.
- 571 Ramsay, J. O., and Ramsey, J. B. (2002) Functional data analysis of the dynamics of the
572 monthly index of nondurable goods production. *Journal of Economics*, **107**, 327–344.
- 573 Tang, Q., and Cheng, L. (2008). M-estimation and B-spline approximation for varying coeffi-
574 cient models with longitudinal data. *Journal of Nonparametric Statistics*, **20**, 611–625.
- 575 Tibshirani, R. (1996). Regression shrinkage and selection vis the Lasso. *Journal of the Royal*

REFERENCES37

- 576 *Statistical Society. Series B (Methodological)*, **58**, 267–288.
- 577 Vogt, M. (2012). Nonparametric regression for locally stationary time series. *The Annals of*
578 *Statistics*, 40(5), 2601–2633.
- 579 Wang, L., and Yang, L. (2007). Spline-backfitted kernel smoothing of nonlinear additive au-
580 toregression model. *The Annals of Statistics*, **35**, 2474–2503.
- 581 Xue, L., Qu, A., and Zhou, J. (2010). Consistent model selection for marginal generalized
582 additive model for correlated data. *Journal of the American Statistical Association*, **105**,
583 1518–1530.
- 584 Xue, L., and Zhu, L. (2007). Empirical likelihood for a varying coefficient model with longitu-
585 dinal data. *Journal of the American Statistical Association*, **102**, 642–654.
- 586 Yao, F. (2007). Asymptotic distributions of nonparametric regression estimators for longitudinal
587 or functional data. *Journal of Multivariate Analysis*, **98**, 40–56.
- 588 Yao, F., Müller, H. G., and Wang, J. L. (2005). Functional linear regression analysis for longi-
589 tudinal data. *Annals of Statistics*, **33**, 2873–2903.
- 590 Zeger, S. L. and Diggle, P. G. (1994). Semiparametric models for longitudinal data with appli-
591 cation to CD4 cell numbers in HIV seroconverters. *Biometrics*, **50**, 689–699.
- 592 Zhang, X., Park, B. U., and Wang J. (2013). Time-varying additive models for longitudinal
593 data. *Journal of the American Statistical Association*, **108**, 983–998.
- 594 Zhang, X., and Wang, J.-L. (2015). Varying-coefficient additive models for functional data.

38REFERENCES

595 *Biometrika*, **102**, 15–32.

596 Zhang, X., and Wang, J.-L. (2016). From sparse to dense functional data and beyond. *The*
597 *Annals of Statistics*, **44**, 2281–2321.

598 Zou, H. (2006). The adaptive lasso and its oracle properties. *Journal of the American Statistical*
599 *Association*, **101**, 1418–1429.

600 Lixia Hu

601 School of Statistics and Mathematics, Shanghai Lixin University of Accounting and Finance,

602 Shanghai, China

603 E-mail: hulx18sufe@163.com

604 Tao Huang

605 School of Statistics and Management, Shanghai University of Finance and Economics, Shanghai,

606 China

607 E-mail: huang.tao@mail.shufe.edu.cn

608 Jinhong You

609 School of Statistics and Management, Shanghai University of Finance and Economics, Shanghai,

610 China

611 E-mail: johnyou07@163.com