# Simultaneous Confidence Corridors and Variable Selection for Generalized Additive Models

Shuzhuan Zheng*
Rong Liu**
Lijian Yang*
Wolfgang Karl Härdle***

BERLIN

ECONOMIC RISK

SFB 649

* Soochow University, China and Michigan State University, USA
** University of Toledo, USA
*** Humboldt-Universität zu Berlin, Germany

# Simultaneous Confidence Corridors and Variable Selection for Generalized Additive Models

Shuzhuan Zheng

Center for Advanced Statistics and Econometrics Research

Soochow University

Suzhou 215006, China

and

Department of Statistics and Probability

Michigan State University

East Lansing, MI 48824 email: `zheng@stt.msu.edu`

Rong Liu

Department of Mathematics and Statistics

University of Toledo

Toledo, OH 43606 email: `rong.liu@utoledo.edu`

Lijian Yang

Center for Advanced Statistics and Econometrics Research

Soochow University

Suzhou 215006, China email: `yanglijian@suda.edu.cn`

and

Department of Statistics and Probability

Michigan State University

East Lansing, MI 48824 email: `yang@stt.msu.edu`

Wolfgang K. Härdle

C.A.S.E. – Center for Applied Statistics and Economics

Humboldt-Universität zu Berlin

Unter den Linden 6

10099 Berlin, Germany email: `haerdle@wiwi.hu-berlin.de`

and

Lee Kong Chian School of Business, Singapore Management University

**Author's Footnote:**

Shuzhuan Zheng is Visiting Scholar, Center for Advanced Statistics and Econometrics Research, Soochow University, Suzhou 215006, China, and Ph.D. student, Department of Statistics and Probability, Michigan State University, East Lansing, MI 48824 (E-mail: `zheng@stt.msu.edu`). Rong Liu is Assistant Professor, Department of Mathematics and Statistics, University of Toledo, Toledo, OH 43606 (E-mail: `rong.liu@utoledo.edu`). Lijian Yang is Director, Center for Advanced Statistics and Econometrics Research, Soochow University, Suzhou 215006, China, and Professor, Department of Statistics and Probability, Michigan State University, East Lansing, MI 48824 (E-mail: `yanglijian@suda.edu.cn;` `yang@stt.msu.edu`). Wolfgang K. Härdle is Professor, C.A.S.E. – Center for Applied Statistics and Economics, Humboldt-Universität zu Berlin, Unter den Linden 6, 10099 Berlin, Germany, and Distinguished Visiting Professor, Lee Kong Chian School of Business, Singapore Management University (E-mail: `haerdle@wiwi.hu-berlin.de`).

## Abstract

In spite of the widespread use of generalized additive models (GAMs), there is no well established methodology for simultaneous inference and variable selection for the components of GAM. There is no doubt that both, inference on the marginal component functions and their selection, are essential in this additive statistical models. To this end, we establish simultaneous confidence corridors (SCCs) and a variable selection criteria through the spline-backfitted kernel smoothing techniques. To characterize the global features of each component, SCCs are constructed for testing their shapes. By extending the BIC to additive models with identity/trivial link, an asymptotically consistent BIC approach for variable selection is proposed. Our procedures are examined in simulations for its theoretical accuracy and performance, and used to forecast the default probability of listed Japanese companies.

4

# 1 INTRODUCTION

The generalized additive model (GAM) has gained popularity on addressing the curse of dimensionality in multivariate nonparametric regressions with non-Gaussian responses. GAM was developed by Hastie and Tibshirani (1990) for blending generalized linear model with nonparametric additive regression, which stipulates that a data set $\left\{ \left( \mathbf{X}_i^\mathsf{T}, Y_i \right) \right\}_{i=1}^n$ consists of iid copies of $\left( \mathbf{X}^\mathsf{T}, Y \right)$ that satisfy:

$$\mathsf{E}(Y|\mathbf{X}) = b' \left\{ m\left(\mathbf{X}\right) \right\}, \mathsf{Var}(Y|\mathbf{X}) = a\left(\phi\right) b'' \left\{ m\left(\mathbf{X}\right) \right\}, m\left(\mathbf{X}\right) = c + \sum_{\alpha=1}^d m_\alpha(X_\alpha), \tag{1}$$

$$Y = b' \left\{ m\left(\mathbf{X}\right) \right\} + \sigma\left(\mathbf{X}\right)\varepsilon, \sigma\left(\mathbf{X}\right) = \left\{ \mathsf{Var}(Y|\mathbf{X}) \right\}^{1/2}$$

where the response $Y$ is one of certain types, such as Bernoulli, Poisson and so forth, the vector $\mathbf{X} = (X_1, X_2, ..., X_d)^\mathsf{T}$ consists of the predictors, $m_\alpha(\cdot), 1 \leq \alpha \leq d$ are unknown smooth functions, the white noise $\varepsilon$ satisfies that $\mathsf{E}\left(\varepsilon|\mathbf{X}\right) = 0$ and $\mathsf{E}\left(\varepsilon^2|\mathbf{X}\right) = 1$, while $c$ is an unknown constant, $a\left(\phi\right)$ is a nuisance parameter that quantifies overdispersion, and $\left(b'\right)^{-1}(\cdot)$ is a known link function. In particular, if one takes the identity/trivial link, model (1) becomes a common additive model, see Huang and Yang (2004).

It is often the case that in model (1) the probability density function of $Y_i$ conditional on $\mathbf{X}_i$ with respect to a fixed $\sigma$-finite measure forms an exponential family:

$$f\left(Y_i|\mathbf{X}_i, \phi\right) = \exp\left[ \left\{ Y_i m\left(\mathbf{X}_i\right) - b\left\{ m\left(\mathbf{X}_i\right) \right\} \right\} / a\left(\phi\right) + h\left(Y_i, \phi\right) \right].$$

Nonetheless, such an assumption is not necessary in this paper. Instead, we only stipulate that the conditional variance and conditional mean are linked by

$$\mathsf{Var}\left(Y|\mathbf{X} = \mathbf{x}\right) = a\left(\phi\right) b'' \left[ \left(b'\right)^{-1} \left\{ \mathsf{E}\left(Y|\mathbf{X} = \mathbf{x}\right) \right\} \right].$$

For identifiability, one needs

$$\mathsf{E}\left\{ m_\alpha\left(X_\alpha\right) \right\} = 0, 1 \leq \alpha \leq d \tag{2}$$

that leads to unique additive representations of $m\left(\mathbf{x}\right) = c + \sum_{\alpha=1}^d m_\alpha\left(x_\alpha\right)$. Without loss of generality, $\mathbf{x}$ take values in $\chi = [0, 1]^d$.

Model (1) has numerous applications. In corporate credit rating, for instance, one is interested in modelling how the default or non-default of a given corporate or company depends on the additive effects of the covariates in financial statements, i.e., the response $Y = 0, 1$ with 1 indicating default, 0 indi-

cating non-default, and the predictors are selected from financial statements with a logit-link $(b')^{-1}(x)$ $= \log\{x/(1-x)\}$. Our method has been applied to $3,472$ companies in Japan within a 5-year default horizon (2005-2010), and it has been discovered that the current liabilities and stock market returns of current, 3 months and 6 months prior to default are very significant as rating factors, and the default impact of the selected factors are examined via the simultaneous confidence corridors (SCCs) in Figure 1 (a)-(c). More details of this example are contained in Section 6.

[Figure 1 about here.]

The smooth functions $\{m_\alpha(x_\alpha)\}_{\alpha=1}^d$ in (1) can be estimated by, for instance, kernel methods in Linton and Härdle (1996), Linton (1997) and Yang, Sperlich and Härdle (2003), B-spline methods in Stone (1986) and Xue and Liang (2010), and two-stage methods in Horowitz and Mammen (2004) and Horowitz et al. (2006). To make statistical inference on these functions individually and collectively, however, the proper tools are simultaneous confidence corridors (SCCs) and consistent variable selection criteria.

The SCC methodology has attracted attention in a variety of applied fields, see Xia (1998), Fan and Zhang (2000), Wu and Zhao (2007), Zhao and Wu (2008), Ma, Yang and Carroll (2012) among others. Capturing shape properties of the functions $\{m_\alpha(x_\alpha)\}_{\alpha=1}^d$ is of utmost importance. A smooth component covered entirely within SCC can be replaced by a parametric one, thereby improving the estimation efficiency, see He, Zhu and Fung (2002), He, Fung and Zhu (2005) for discussions. To our knowledge, SCCs have not been established due to a technological lack of estimators that fit in Gaussian process extreme value theory. Using the spline-backfitted kernel (SBK) methodology of Liu, Yang and Härdle (2013) (hereafter LYH), we extend work of univariate nonparametric regression in Bickel and Rosenblatt (1973) and Härdle (1989) to those of GAM. The SBK technique has been studied in Wang and Yang (2007), Wang and Yang (2009), Liu and Yang (2010) and Ma and Yang (2011) for the simpler additive model (i.e., GAM with $b'(x) \equiv x$) including the construction of SCC, but ours is the first work on SCC for GAM with nonlinear link.

While variable selection for nonparametric additive model has been investigated under different settings, see Wang, Li and Huang (2008), there is lack of theoretically reliable variable selection for GAM. To the best of our knowledge, only Zhang and Lin (2006) proposed the "COSSO" method for variable selection in nonparametric regression with exponential families, but without asymptotic theory. Instead, we tackle this issue by building a BIC type criterion based on spline pre-smoothing (first stage in the SBK), which is asymptotically consistent and easy to compute. Our work extends the BIC criterion for additive models

(trivial link) in Huang and Yang (2004). This extension is challenging since a much more complicated quasi-likelihood is employed with nonlinear link instead of the log mean squared error for trivial link. The appendix gives more details.

The rest of paper is organized as follows. The SBK estimator and its oracle property are briefly described in Section 2. Asymptotic extreme value distribution of the SBK estimator is investigated in Section 3, which is used to construct the SCCs of component functions. Section 4 introduces a BIC criterion in the GAM setting and provides results on consistent component selection as well as the implementation, followed by the Monte Carlo simulations in Section 5. Section 6 illustrates the application of our SCC and BIC methods to predict default of nearly $3,500$ listed companies in Japan. Technical assumptions and proofs are presented in the Appendix.

## 2 SPLINE-BACKFITTED KERNEL SMOOTHING IN GAM

In this section we briefly describe the SBK estimator for GAM (1) and its oracle properties obtained in LYH. Let $\{\mathbf{X}_i, Y_i\}_{i=1}^n$ be i.i.d. observations following model (1). Without loss of generality, one denotes $\mathbf{x}_{-1} = (x_2, ..., x_d)$ and $m_{-1}(\mathbf{x}_{-1}) = c + \sum_{\alpha=2}^{d} m_\alpha(x_\alpha)$ and estimates $m_1(x_1)$.

As a benchmark of efficiency, we introduce the "oracle smoother" by treating the constant $c$ and the last $d-1$ components $\{m_\alpha(x_\alpha)\}_{\alpha=2}^d$ as known. The only unknown component $m_1(x_1)$ is estimated by maximizing a local log-likelihood function $\widetilde{l}(a, x_1)$ for each $x_1 \in [h, 1-h]$:

$$\widetilde{l}(a, x_1) = n^{-1} \sum_{i=1}^{n} [Y_i \{a + m_{-1}(\mathbf{X}_{i,-1})\} - b\{a + m_{-1}(\mathbf{X}_{i,-1})\}] K_h(X_{i1} - x_1), \qquad (3)$$

where $a \in A$, a set whose interior contains $m_1([0,1])$. The oracle smoother of $m_1(x_1)$ is

$$\widetilde{m}_{\mathrm{K},1}(x_1) = \underset{a \in A}{\operatorname{argmax}} \, \widetilde{l}(a, x_1). \qquad (4)$$

Although $\widetilde{m}_{\mathrm{K},1}(x_1)$ is not a statistic since $c$ and $\{m_\alpha(x_\alpha)\}_{\alpha=2}^d$ are actually unknown, its asymptotic properties serve as a benchmark for estimators of $m_1(x_1)$ to achieve.

To define the SBK, we introduce the linear B spline basis for smoothing: $b_J(x) = (1 - |x - \xi_J|/H)_+$, $0 \le J \le N+1$ where $0 = \xi_0 < \xi_1 < \cdots < \xi_N < \xi_{N+1} = 1$ are a sequence of equally spaced points, called interior knots, on interval $[0,1]$. Denote by $H = (N+1)^{-1}$ the width of each subinterval $[\xi_J, \xi_{J+1}], 0 \le J \le N$ and the degenerate knots by $\xi_{-1} = 0, \xi_{N+2} = 1$. The space of $\alpha$-empirically centered linear spline

functions on $[0, 1]$ is

$$G_{n,\alpha}^0 = \left\{ g_\alpha : g_\alpha(x_\alpha) = \sum_{J=0}^{N+1} \lambda_J b_J(x_\alpha), \mathsf{E}_n\{g_\alpha(X_\alpha)\} = 0 \right\}, 1 \leq \alpha \leq d, \tag{5}$$

with empirical expectation $\mathsf{E}_n\{g_\alpha(X_\alpha)\} = n^{-1}\sum_{i=1}^n g_\alpha(X_{\alpha i})$. The space of additive spline functions on $\chi = [0,1]^d$ is

$$G_n^0 = \left\{ g(\mathbf{x}) = c + \sum_{\alpha=1}^d g_\alpha(x_\alpha); c \in \mathbb{R}, \ g_\alpha \in G_{n,\alpha}^0 \right\}. \tag{6}$$

The SBK method is defined in two steps. One first pre-estimates the unknown functions $\{m_\alpha(x_\alpha)\}_{\alpha=2}^d$ and constants $c$ by linear spline smoothing. We define the log-likelihood function $\widehat{L}(g)$ as

$$\widehat{L}(g) = n^{-1}\sum_{i=1}^n [Y_i g(\mathbf{X}_i) - b\{g(\mathbf{X}_i)\}], g \in G_n^0. \tag{7}$$

According to Lemma 14 of Stone (1986), (7) has a unique maximizer with probability approaching 1. Therefore, the multivariate function $m(\mathbf{x})$ can be estimated by an additive spline function:

$$\widehat{m}(\mathbf{x}) = \underset{g \in G_n^0}{\operatorname{argmax}}\, \widehat{L}(g). \tag{8}$$

The spline estimator is asymptotically consistent, and can be calculated efficiently. However, no measure of confidence can be assigned to the spline estimator, see Wang and Yang (2007) and LYH. To overcome this problem, we adapt the SBK estimator, which combines the strength of kernel smoothing with regression spline. One then rewrites $\widehat{m}(\mathbf{x}) = \widehat{c} + \sum_{\alpha=1}^d \widehat{m}_\alpha(X_{i\alpha})$ for $\widehat{c} \in \mathbb{R}$ and $\widehat{m}_\alpha(x_\alpha) \in G_{n,\alpha}^0$ and defines a univariate quasi-likelihood function similar to $\widetilde{l}(a, x_1)$ in (3) as

$$\widehat{l}(a, x_1) = n^{-1}\sum_{i=1}^n [Y_i\{a + \widehat{m}_{\text{-}1}(\mathbf{X}_{i,\text{-}1})\} - b\{a + \widehat{m}_{\text{-}1}(\mathbf{X}_{i,\text{-}1})\}] K_h(X_{i1} - x_1), \tag{9}$$

with $\widehat{m}_{\text{-}1}(\mathbf{x}_{\text{-}1}) = \widehat{c} + \sum_{\alpha=2}^d \widehat{m}_\alpha(x_\alpha)$ being the pilot spline estimator of $m_{\text{-}1}(\mathbf{x}_{\text{-}1})$. Consequently, the SBK estimator of $m_1(x_1)$ is

$$\widehat{m}_{\text{SBK},1}(x_1) = \underset{a \in A}{\operatorname{argmax}}\, \widehat{l}(a, x_1). \tag{10}$$

We now introduce some useful results and definitions from LYH, under Assumptions (A1)-(A7) in appendix, as $n \to \infty$,

$$\sup_{x_1 \in [0,1]} |\widehat{m}_{\text{SBK},1}(x_1) - \widetilde{m}_{\text{K},1}(x_1)| = \mathcal{O}_{a.s.}\left(n^{-1/2}\log n\right), \tag{11}$$

8

$$\widetilde{m}_{\mathrm{K},1}\left(x_1\right) - m_1\left(x_1\right) = \mathrm{bias}_1\left(x_1\right) h^2/D_1\left(x_1\right) + n^{-1}\sum_{i=1}^{n} K_h\left(X_{i1} - x_1\right)\sigma\left(X_i\right)\varepsilon_i/D_1\left(x_1\right) + r_{\mathrm{K},1}\left(x_1\right) \quad (12)$$

in which the higher order remainder $r_{\mathrm{K},1}\left(x_1\right)$ satisfies

$$\sup_{x_1\in[h,1-h]}\left|r_{\mathrm{K},1}\left(x_1\right)\right| = \mathcal{O}_{a.s.}\left(n^{-1/2}h^{1/2}\log n\right). \quad (13)$$

The scale function $D_1\left(x_1\right)$ and bias function $\mathrm{bias}_1\left(x_1\right)$ are defined in LYH as:

$$\sigma_b^2\left(x_1\right) = \mathsf{E}\left[b''\left\{m\left(\mathbf{X}\right)\right\}|X_1 = x_1\right], \ \sigma^2\left(x_1\right) = \mathsf{E}\left\{\sigma^2\left(\mathbf{X}\right)|X_1 = x_1\right\} \quad (14)$$

$$D_1\left(x_1\right) = f_1\left(x_1\right)\sigma_b^2\left(x_1\right), v_1^2\left(x_1\right) = \|K\|_2^2 f_1\left(x_1\right)\sigma^2\left(x_1\right). \quad (15)$$

$$\mathrm{bias}_1\left(x_1\right) = \mu_2\left(K\right)\times \quad (16)$$
$$\left\{m_1''\left(x_1\right)D_1\left(x_1\right) + m_1'\left(x_1\right)f\left(x_1\right)\sigma_b^2\left(x_1\right)' - \left\{m_1'\left(x_1\right)\right\}^2 f\left(x_1\right)\mathsf{E}\left[b'''\left\{m\left(\mathbf{X}\right)\right\}|X_1 = x_1\right]\right\}$$

where $\|K\|_2^2 = \int K^2\left(u\right)du$, $\mu_2\left(K\right) = \int K\left(u\right)u^2 du$. The above equations (11), (12) and (13) lead one to a simplifying decomposition of the estimation error $\widehat{m}_{\mathrm{SBK},1}\left(x_1\right) - m_1\left(x_1\right)$

$$\sup_{x_1\in[h,1-h]}\left|\widehat{m}_{\mathrm{SBK},1}\left(x_1\right) - m_1\left(x_1\right) - n^{-1}\sum_{i=1}^{n}K_h\left(X_{i1} - x_1\right)\sigma\left(\mathbf{X}_i\right)\varepsilon_i/D_1\left(x_1\right)\right| \quad (17)$$
$$= \mathcal{O}_{a.s.}\left(n^{-1/2}h^{1/2}\log n + n^{-1/2}\log n + h^2\right).$$

A decomposition such as (17) has not appeared in the literature for any other estimators of $m_1\left(x_1\right)$, and it is fundamental for constructing SCCs in section 3.

## 3   GAM INFERENCE VIA SCC

In this section, we propose SCCs for GAM components.

### 3.1   Main Results

Denote $a_h = \sqrt{-2\log h}$, $C\left(K\right) = \|K'\|_2^2\|K\|_2^{-2}$ and for any $\alpha \in \left(0,1\right)$, the quantile

$$Q_h(\alpha) = a_h + a_h^{-1}\left[\log\left\{\sqrt{C\left(K\right)}/\left(2\pi\right)\right\} - \log\left\{-\log\sqrt{1-\alpha}\right\}.\right] \quad (18)$$

Also with $D_1(x_1)$ and $v_1^2(x_1)$ given in (15), we define

$$\sigma_n(x_1) = n^{-1/2}h^{-1/2}v_1(x_1)D_1^{-1}(x_1). \tag{19}$$

THEOREM 1  *Under Assumptions (A1)-(A7), as $n \to \infty$*

$$\lim_{n\to\infty} \mathrm{P}\left\{\sup_{x_1\in[h,1-h]}|\widehat{m}_{\mathrm{SBK},1}(x_1) - m_1(x_1)|/\sigma_n(x_1) \le Q_h(\alpha)\right\} = 1 - \alpha.$$

*A $100(1-\alpha)\%$ simultaneous confidence corridor for $m_1(x_1)$,*

$$\widehat{m}_{\mathrm{SBK},1}(x_1) \pm \sigma_n(x_1)Q_h(\alpha). \tag{20}$$

The above SCC for component function $m_1(x_1)$ resembles the SCCs in Bickel and Rosenblatt (1973) and Härdle (1989) for estimating unknown univariate nonparametric function, although it is for high dimensional nonparametric regression.

## 3.2  Implementation

To construct the SCC for $m_1(x_1)$ in (20), one needs to select the bandwidth $h$ first, and then evaluate $m_{\mathrm{SBK},1}(x_1)$, $Q_h(\alpha)$ and $\sigma_n(x_1)$ given in (10), (18) and (19).

Assumption (A6) requires that the bandwidth for SCC be slightly smaller than the mean square optimal bandwidth $h_{\mathrm{opt}}$ (minimizing AMISE) in LYH, we therefore have taken $h = h_{\mathrm{opt}}(\log n)^{-1/4}$ as a data-driven undersmoothing bandwidth for SCC construction to fulfill Assumption (A6). Recent articles on SCC for time series, such as Wu and Zhao (2007), Zhao and Wu (2008), have used similar undersmoothing bandwidths.

For a given $\alpha$ and a chosen bandwidth $h$, one can easily estimate $m_{\mathrm{SBK},1}(x_1)$ and $Q_h(\alpha)$ as in (10), (18). To evaluate $\sigma_n(x_1)$, one needs to estimate $v_1(x_1)$ and $D_1^{-1}(x_1)$ given in (15), i.e., estimating $f(x_1), \sigma_b^2(x_1)$ and $\sigma^2(x_1)$. The estimation of the density function $f(x_1)$ is trivial, namely, $\widehat{f}(x_1) = n^{-1}\sum_{i=1}^n K_h(X_{i1} - x_1)$. We further illustrate the spline estimates of $\sigma_b^2(x_1)$ and $\sigma^2(x_1)$ below:

One partitions $\min_i X_{i1} = t_{1,0} < \cdots < t_{1,N+1} = \max_i X_{i1}$ where $N$ is the number of spline interior knots, i.e., $N = N_n = \min\left([n^{1/4}\log n] + 1, [n/4d - 1/d] - 1\right)$ which satisfies the assumption (A7) in the Appendix. Then $\sigma_b^2(x_1)$ can be estimated as $\sum_{k=0}^3 \widehat{a}_{1,k}^k x_1^k + \sum_{k=4}^{N+3} \widehat{a}_{1,k}(x_1 - t_{\alpha,k-3})^3$ where $\{\widehat{a}_{1,k}\}_{k=0}^{N+3}$

minimize

$$\sum_{i=1}^{n} \left[ b'' \left\{ \widehat{m} \left( \mathbf{X}_i \right) \right\} - \left\{ \sum_{k=0}^{3} a_{1,k} X_{i1}^k + \sum_{k=4}^{N+3} a_{1,k} \left( X_{i1} - t_{k-3} \right)^3 \right\} \right]^2, \tag{21}$$

and $\sigma^2(x_1)$ can be estimated as $\sum_{k=0}^{3} \widehat{a}_{1,k}^k x_1^k + \sum_{k=4}^{N+3} \widehat{a}_{1,k} \left( x_1 - t_{\alpha,k-3} \right)^3$ where $\{\widehat{a}_{1,k}\}_{k=0}^{N+3}$ minimize

$$\sum_{i=1}^{n} \left[ \left[ Y_i - b' \left\{ \widehat{m} \left( \mathbf{X}_i \right) \right\} \right]^2 - \left\{ \sum_{k=0}^{3} a_{\alpha,k} X_{i1}^k + \sum_{k=4}^{N+3} a_{\alpha,k} \left( X_{i1} - t_{k-3} \right)^3 \right\} \right]^2. \tag{22}$$

The resulted estimate $\hat{\sigma}_n(x_1)$ of $\sigma_n(x_1)$, using (21) and (22) satisfies $\sup_{x_1 \in [h,1-h]} |\hat{\sigma}_n(x_1) - \sigma_n(x_1)| = \mathcal{O}_p(n^{-\gamma})$ for some $\gamma > 0$, see LYH Section 5 for details. This consistency and Slutzky's theorem ensure that $\mathrm{P} \left\{ \sup_{x_1 \in [h,1-h]} |\widehat{m}_{\mathrm{SBK},1}(x_1) - m_1(x_1)| / \hat{\sigma}_n(x_1) \leq Q_h(\alpha) \right\} \to 1 - \alpha$ as $n \to \infty$, and therefore

$$\widehat{m}_{\mathrm{SBK},1}(x_1) \pm \hat{\sigma}_n(x_1) Q_h(\alpha) \tag{23}$$

is a $100(1-\alpha)\%$ simultaneous confidence corridor for $m_1(x_1)$. The SCC constructions of other components $m_2(x_2), ..., m_d(x_d)$ are similar. It is worthwhile to emphasize that, in general, the estimators $\widehat{m}_{\mathrm{SBK},1}(x_1)$, $\widehat{Q}_h(\alpha), \widehat{f}(x_1)$ and $\hat{\sigma}_n(x_1)$ remain stable if $h$ slightly varies.

# 4   VARIABLE SELECTION IN GAM

In this section, we propose a Bayesian Information Criterion (BIC) for component function selection based on spline smoothing for GAM and an efficient implementation follows.

## 4.1   Main Results

According to Stone (1985), p.693, the space of $\alpha$-centered square integrable functions on $[0,1]$ is defined as

$$\mathcal{H}^0 = \left\{ g : \mathsf{E} \left\{ g \left( X_\alpha \right) \right\} = 0, \mathsf{E} \left\{ g^2 \left( X_\alpha \right) \right\} < \infty, 1 \leq \alpha \leq d \right\}, \tag{24}$$

and the model space $\mathcal{M}$ is

$$\mathcal{M} = \left\{ g(\mathbf{x}) = c + \sum_{\alpha=1}^{d} g_\alpha \left( \mathbf{x}_\alpha \right); g_\alpha \in \mathcal{H}^0, 1 \leq \alpha \leq d \right\}, \tag{25}$$

where $c$ is a finite constant.

To introduce the proposed BIC, let $\{1, \ldots, d\}$ denote the complete set of indices of $d$ tuning variables $(X_1, \ldots, X_d)$. For each subset $S \subset \{1, \ldots, d\}$, define a corresponding model space $\mathcal{M}_S$ for $S$ as

$$\mathcal{M}_S = \left\{ g(\mathbf{x}) = c + \sum_{\alpha \in S} g_\alpha(\mathbf{x}_\alpha) ; c \in \mathbb{R}, g_\alpha \in \mathcal{H}^0, \alpha \in S \right\}, \tag{26}$$

with $\mathcal{H}^0$ given in (24), and the space of the additive spline functions as

$$G_{n,S}^0 = \left\{ g(\mathbf{x}) = c + \sum_{\alpha \in S} g_\alpha(x_\alpha) ; c \in \mathbb{R}, g_\alpha \in G_{n,\alpha}^0, \alpha \in S \right\}, \tag{27}$$

with $G_{n,\alpha}^0$ given in (5). Define the least squares projection of function $m$ in $\mathcal{M}_S$ as

$$m_S = \underset{g \in \mathcal{M}_S}{\arg\min} \, \mathsf{E} \left\{ m(\mathbf{X}) - g(\mathbf{X}) \right\}^2 \tag{28}$$

and define the set $S_0$ of significant variables as the minimal set $S \subset \{1, \ldots, d\}$ such that $\mathsf{E}\{m(\mathbf{X}) - m_S(\mathbf{X})\}^2 = 0$, which is uniquely defined according to Lemma 1 of Huang and Yang (2004).

To identify $S_0$, one computes for an index set $S$ the BIC as

$$\mathrm{BIC}_S = -2\widehat{L}(\widehat{m}_S) + \frac{N_S}{n} (\log n)^3 \tag{29}$$

where $\widehat{L}(\cdot)$ is given in (7), $\widehat{m}_S(\mathbf{x}) \in G_{n,S}^0$ is the pilot spline estimator as in (8), $N_S = 1 + (N+1) \#(S)$ with $N$ the number of interior knots, $\#(S)$ the cardinality of $S$. In practice, $N = N_n$ can be taken as

$$\min\left( \left[ n^{1/4} \log n \right] + 1, [n/4d - 1/d] - 1 \right), \tag{30}$$

which satisfies the assumption (A7) in the Appendix.

Our variable selection rule takes the subset $\widehat{S} \subset \{1, \ldots, d\}$ that minimizes $\mathrm{BIC}_S$.

THEOREM **2** *Under Assumptions (A1)-(A5) and (A7),* $\lim_{n \to \infty} \mathrm{P}\left( \widehat{S} = S_0 \right) = 1.$

According to Theorem 2, the variable selection rule based on the BIC in (29) is consistent. The nonparametric version BIC was firstly established in Huang and Yang (2004) for additive autoregression model, and adapted to additive coefficient model by Xue and Yang (2006), to single index model by Wang and Yang (2009). Our proposed BIC differs from all of the above as it is based on quasi-likelihood rather than mean squared error, which makes the technical proof of consistency much more challenging. To the

best of our knowledge, it is the first theoretically reliable information criterion in this setting.

## 4.2 Implementation

The proposed BIC is implemented without a greedy search through all possible subsets. Instead, the forward stepwise regression procedure is used with minimizing BIC as the criterion.

# 5 MONTE CARLO SIMULATION

This section studies the performance of the proposed procedures, reporting also the computational costs, the consistency of selecting variables via BIC and the global coverage precision of the SCC. The data are generated from

$$\mathrm{P}(Y = 1 | \mathbf{X} = \mathbf{x}) = b' \left\{ c + \sum_{\alpha=1}^{d} m_\alpha (X_\alpha) \right\}, b'(x) = \frac{e^x}{1 + e^x} \tag{31}$$

with $d = 10, c = 0, m_3(x) = m_4(x) = m_5(x) = \sin(\pi x)$, $m_6(x) = x, m_7(x) = e^x - (e - e^{-1})$ and $m_\alpha(x) = 0$ for $\alpha = 1, 2, 8, 9, 10$. The set $S_0$ is therefore $S_0 = \{3, 4, 5, 6, 7\}$.

The predictors are generated by

$$X_{i\alpha} = 2\Phi(Z_{i\alpha}) - 1, \ \mathbf{Z}_i = (Z_{i1}, ..., Z_{id}) \sim \mathrm{N}(0, \Sigma), 1 \le i \le n, 1 \le \alpha \le d, \tag{32}$$

where $\Phi$ is the standard normal c.d.f. and $\Sigma = (1 - r) \mathbf{I}_{d \times d} + r \mathbf{1}_d \mathbf{1}_d^T$. The parameter $r$ $(0 \le r < 1)$ controls the correlation between $Z_{i\alpha}, 1 \le \alpha \le d$.

In what follows, the performance of BIC and COSSO is firstly compared, followed by a computational comparison between the SBK and a kernel method in GAM, and it ends with a report on the SCC global coverage for components.

[Table 1 about here.]

[Table 2 about here.]

Table 1 shows the simulation results from 100 replications, where the outcome is defined in accuracy as correct fitting, if $\widehat{S} = S_0$; overfitting, if $S_0 \subset \widehat{S}$; and underfitting, if $S_0 \nsubseteq \widehat{S}$. It is clear that the performance of BIC on selecting 5 significant variables $m_\alpha(X_\alpha), \alpha = 3, ..., 7$, is quite satisfactory. The selection accuracy becomes higher as the sample size increases and the correlation decreases. The accuracy and computing

time of COSSO are also listed for comparison (Platform: R; PC: Intel 3.4 GHz processor and 16 GB RAM). It is obvious that the BIC significantly outperforms the COSSO in terms of both accuracy and computing time. To examine the computing advantage of BIC for large $d$, we extend $d = 10$ to $50$ by using $m_3, ..., m_7$ as above and all the other component functions are 0. The BIC is vastly faster than COSSO for $d = 50$. All of these findings confirm what is expected according to the asymptotic theory.

The SCC global coverage for $m_\alpha(x_\alpha), \alpha = 3, ..., 7$ is reported in Table 2. It turns out that the empirical coverage approaches the nominal confidence levels as $n$ increases, and a better coverage occurs when the correlation is lower.

The above studies evidently indicate the reliability of our methodology, such as a high selection accuracy of the BIC and a desired global coverage of the SCC, which make their applications for credit rating modelling in the following section sensible.

# 6 APPLICATION

We now return to forecast default probabilities of the listed companies in Japan. The data provided by the Risk Management Institute, National University of Singapore include the comprehensive financial statements and the credit events (default or bankruptcy) from 2005 to 2010 of $3,583$ Japanese firms.

Berg (2007) found that the liability status was important to indicate the creditworthiness of a company, while Bernhardsen (2001) and Ryser and Denzler (2009) proposed to consider the "leverage effect" expressed by the financial statement ratios. Therefore, we have pooled two situations by considering $X_1$: Current liability, $X_2$: Current stock return, $X_3$: Long term borrow, $X_4$: Short term borrow, $X_5$: Total asset, $X_6$: Non-current liability, $X_7$: 3 months earlier (stock) return, $X_8$: 6 months earlier (stock) return, $X_9$: Current ratio, $X_{10}$: Net liability to shareholder equity, $X_{11}$: Shareholder equity to total liability and equity, $X_{12}$: TCE ratio, $X_{13}$: Total debt to total asset, $X_{14}$: Quick ratio.

Selecting the rating factors via the BIC given in (29), we have found that $X_1$: Current liabilities, $X_7$: 3 months earlier return, $X_8$: 6 months earlier return are selected. Similar rating covariates were also discovered in Shina and Moore (2003), Berg (2007) and Ryser and Denzler (2009). However, Berg (2007) selected 23 variables which led to a non-parsimonious GAM. In contrast, Ryser and Denzler (2009) had found that 3 financial ratios (capital turnover, long-term debt ratio, return on total capital) were significant based on the blockwise cross-validation (CV) method which is nonetheless extremely time consuming in comparison to the proposed BIC.

Figure 1 (a)-(c) depicts the SBK estimator of the factor's default impact curve on domain, while a shoal of 95% CIs and the 95% SCC present respectively the pointwise and global uncertainty of the whole curve. The SBK estimators indicate overall monotonicities of each rating factors, and the SCCs turn out to be fairly narrow to warrant the global nonlinearities of the factors' curves which reveal the underlying nonlinear features in different segments of domain.

As for the model evaluations, the Cumulative Accuracy Profile (CAP) is used. For any credit scoring method $S$, one defines its alarm rate $F(s) = \mathrm{P}(S \leq s)$ and the hit rate $F_{\mathrm{D}}(s) = \mathrm{P}(S \leq s \,|\, \mathrm{D})$ where $\mathrm{D}$ represents the conditioning event of "default". One then defines the CAP curve for $S$ as

$$\mathrm{CAP}(u) = F_{\mathrm{D}}\left\{F^{-1}(u)\right\}, u \in (0,1), \tag{33}$$

which is the percentage of default-infected obligators that are found among the first (according to their scores) $100u\%$ of all obligators. The perfect scoring method $\mathrm{P}$ assigns low scores first to all default-infected obligators and thus $\mathrm{CAP}_{\mathrm{P}}(u) = \min(u/p, 1), u \in (0,1)$ where $p$ is the unconditional default probability, whereas the completely noninformative scoring method with zero discriminatory power displays a diagonal line $\mathrm{CAP}_{\mathrm{N}}(u) \equiv u, u \in (0,1)$, see details of the CAP in Engelmann, Hayden and Tasche (2003).

A satisfactory scoring method's CAP curve would be expected to approach $\mathrm{CAP}_{\mathrm{P}}(u)$ and always better than the noinformative $\mathrm{CAP}_{\mathrm{N}}(u)$, and one uses the Accuracy Ratio (AR) to quantify its position. The AR is the ratio of the area $a_R$ enclosed between the given CAP curve and the noninformative diagonal curve $\mathrm{CAP}_{\mathrm{N}}(u) \equiv u$, and the total area $a_P$ enclosed between the perfect CAP curve $\mathrm{CAP}_{\mathrm{P}}(u)$ and the noninformative diagonal curve $\mathrm{CAP}_{\mathrm{N}}(u)$. Thus

$$\mathrm{AR} = \frac{a_R}{a_P} = \frac{2 \int_0^1 \mathrm{CAP}(u)\, du - 1}{1 - p},$$

where $\mathrm{CAP}(u)$ is given in (33). The AR takes value in $[0, 1]$, with value 0 corresponding to the noninformative scoring, and 1 the perfect scoring method, a higher AR indicates an overall higher discriminatory power of a method.

We have applied both GAM and GLM to the first 2000 companies and computed default probabilities of the remaining 1583 companies, and used the default probabilities as scores. Figure 1 (d) depicts the CAPs of GAM (thick solid) and GLM (thin solid), as well as the $\mathrm{CAP}_{\mathrm{P}}(u)$ (dashed) and $\mathrm{CAP}_{\mathrm{N}}(u)$ (dotted). Visually it is clear that GAM has much higher accuracy ratio than GLM, which is in fact the case: the AR is 97.56% for GAM, much higher than the 89.76% for GLM. Nonetheless, we failed to apply the COSSO

for the same data.

# APPENDIX

In what follows, we take $\|\cdot\|$ and $\|\cdot\|_\infty$ as the Euclidean and supremum norms, respectively, i.e., for any $\mathbf{x} = (x_1, x_2, ..., x_d)^{\mathsf{T}} \in \mathbb{R}^d, \|\mathbf{x}\| = \left(\sum_{\alpha=1}^d x_\alpha^2\right)^{1/2}$ and $\|\mathbf{x}\|_\infty = \max_{1 \le \alpha \le d} |x_\alpha|$. For any interval $[a, b]$, denote the space of $p$-th order smooth function by $C^{(p)}[a, b] = \{g \, | g^{(p)} \in C[a, b]\}$, and the class of Lipschitz continuous functions by $\text{Lip}([a, b], C) = \{g \, | |g(x) - g(x')| \le C|x - x'|, \forall x, x' \in [a, b]\}$ for constant $C > 0$. Lastly, define the following latent regression errors

$$\xi_i = Y_i - b'\{m(\mathbf{X}_i)\} = \sigma(\mathbf{X}_i)\varepsilon_i, 1 \le i \le n. \tag{A.1}$$

We need the following technical assumptions:

(A1) *The additive component functions* $m_\alpha \in C^{(1)}[0,1], 1 \le \alpha \le d$: $m_1 \in C^{(2)}[0,1], m'_\alpha \in \text{Lip}([0,1], C_m)$, $2 \le \alpha \le d$ *for some constant* $C_m > 0$.

(A2) *The inverse link function* $b'$ *satisfies that* $b' \in C^2(\mathbb{R}), b''(\theta) > 0, \theta \in \mathbb{R}$. *For a compact interval* $\Theta$ *whose interior contains* $m\left([0,1]^d\right)$, $C_b > \max_{\theta \in \Theta} b''(\theta) \ge \min_{\theta \in \Theta} b''(\theta) > c_b$ *for constants* $0 < c_b < C_b < \infty$.

(A3) *The conditional variance function* $\sigma^2(\mathbf{x})$ *is continuous and positive for* $\mathbf{x} \in [0,1]^d$. *The errors* $\{\varepsilon_i\}_{i=1}^n$ *satisfy that* $\mathsf{E}(\varepsilon_i | \mathbf{X}_i) = 0, \mathsf{E}\left(|\varepsilon_i|^{2+\eta}\right) \le C_\eta$ *for some* $\eta \in (1/2, 1]$.

(A4) *The joint density* $f(\mathbf{x})$ *of* $(X_1, ..., X_d)$ *is continuous:* $0 < c_f \le \inf_{\mathbf{x} \in [0,1]^d} f(\mathbf{x}) \le \sup_{\mathbf{x} \in [0,1]^d} f(\mathbf{x}) \le C_f < \infty$. *The marginal density function* $f_\alpha(x_\alpha)$ *of* $X_\alpha$ *have continuous derivatives on* $[0,1]$ *and the uniform bounds* $C_f$ *and* $c_f$. *There exists a* $\sigma$-*finite measure* $\lambda$ *on* $\mathbb{R}$ *such that the distribution of* $Y_i$ *conditional on* $\mathbf{X}_i$ *has a probability density function* $f_{Y|\mathbf{X}}(y; b'\{m(\mathbf{x})\})$ *relative to* $\lambda$ *whose support for* $y$ *is a common* $\Omega$, *and is uniformly continuous in* $\mathbf{x} \in [0,1]^d$ *for* $y \in \Omega$.

(A5) $\{\mathbf{Z}_i = (\mathbf{X}_i^{\mathsf{T}}, \varepsilon_i)\}_{i=1}^n$ *are independent and identically distributed.*

(A6) *The kernel function* $K(x)$ *is a symmetric probability density function supported on* $[-1, 1]$ *and* $\in C^1[-1, 1]$. *The bandwidth* $h = h_n$ *satisfies that* $h = o\{n^{-1/5}(\log n)^{-1/5}\}, h^{-1} = \mathcal{O}\{n^{1/5}(\log n)^\delta\}$ *for some constant* $\delta > 1/5$.

(A7) *The number of interior knots $N$ satisfies that $c_N n^{1/4} \log n \leq N \leq C_N n^{1/4} \log n$ for some constants $c_N, C_N > 0$.*

Assumptions (A1)-(A7) are standard in GAM, see Stone (1986), Xue and Liang (2010). Assumptions (A5), (A6) are more restrictive than in LYH for the purpose of constructing SCCs, but are unnecessary for Theorem 2 on the consistency of BIC.

## A.1. Preliminaries

Throughout this section, $C$ denotes some generic positive constant unless stated otherwise. Define

$$M_h(t) = h^{-1/2} \int_0^1 K\{(x - t)/h\} \, dW(x) \tag{A.2}$$

where $W(x)$ is a Wiener process defined on $(0, \infty)$ and denote

$$d_h = (-2 \log h)^{1/2} + (-2 \log h)^{-1/2} \left\{ \sqrt{C(K)}/(2\pi) \right\} \tag{A.3}$$

with $C(K)$ given in (18).

LEMMA A.1 *Under Assumption (A6). for any $x \in \mathbb{R}$*

$$\lim_{n \to \infty} \mathrm{P}\left[ (-2 \log h)^{1/2} \left\{ \sup_{t \in [h, 1-h]} |M_h(t)| / \|K\|_2^2 - d_h \right\} < x \right] = e^{-2e^{-x}}.$$

PROOF. One simply applies the same steps in proving Lemma 2.2 of Härdle (1989).

Denote by $T_i$ the random variable $b'\{m(\mathbf{X}_i)\}$, and the Lebesgue measure on $\mathbb{R}^d$ as $\mu^{(d)}$. By Assumption (A4), $\mathbf{X}_i$ has pdf wrt the Lebesgue measure $\mu^{(d)}$, and Assumptions (A1) and (A2) ensure that functions $b'$ and $m$ are at least $C^1$, thus the random vector $(T_i, X_{i1})$ has a joint pdf wrt the Lebesgue measure $\mu^{(2)}$, which one denotes as $f_{T, X_1}(t, x_1)$.

LEMMA A.2 *Under Assumptions (A1)-(A5), for $\xi_i$ in (A.1), the distribution of $(\xi_i, X_{i1})$ has joint pdf wrt $\mu^{(2)}$ as $f_{\xi, X_1}(z, x_1) = \int_\Omega f_{Y|\mathbf{X}}(y; y - z) f_{T, X_1}(y - z, x_1) \, d\lambda(y)$.*

PROOF. The joint pdf of $(Y_i, T_i, X_{i1})$ wrt $\lambda \times \mu^{(2)}$ is therefore $f_{Y|\mathbf{X}}(y; t) f_{T, X_1}(t, x_1)$. For any $(z, x_1) \in \mathbb{R} \times [0, 1]$, and $\triangle z, \triangle x_1 > 0$, one has

$$\mathrm{P}\left[ (\xi_i, X_{i1}) \in (z - \triangle z, z + \triangle z) \times (x_1 - \triangle x_1, x_1 + \triangle x_1) \right] =$$

17

$$P\left[(Y_i - T_i, X_{i1}) \in (z - \triangle z, z + \triangle z) \times (x_1 - \triangle x_1, x_1 + \triangle x_1)\right] =$$

$$= \int_\Omega d\lambda(y) \int_{y - \tau \in (z - \triangle z, z + \triangle z)} d\tau \int_{\chi_1 \in (x_1 - \triangle x_1, x_1 + \triangle x_1)} f_{Y|\mathbf{X}}(y; \tau) f_{T, X_1}(\tau, \chi_1) d\chi_1.$$

Applying dominated convergence theorem, one has as $\max(\triangle z, \triangle x_1) \to 0$

$$\left| P\left[(\xi_i, X_{i1}) \in (z - \triangle z, z + \triangle z) \times (x_1 - \triangle x_1, x_1 + \triangle x_1)\right] - \left\{ \int_\Omega f_{Y|\mathbf{X}}(y; y - z) f_{T, X_1}(y - z, x_1) d\lambda(y) \right\} \right.$$

$$\left. \times \mu^{(2)}\left[(z - \triangle z, z + \triangle z) \times \{(x_1 - \triangle x_1, x_1 + \triangle x_1) \cap [0, 1]\}\right] \right| = o(1)$$

hence the the joint pdf of $(\xi_i, X_{i1})$ wrt $\mu^{(2)}$ is $\int_\Omega f_{Y|\mathbf{X}}(y; y - z) f_{T, X_1}(y - z, x_1) d\lambda(y)$.

For theoretical analysis, we write $c_{J,\alpha} = \mathsf{E}\, b_J(X_\alpha) = \int b_J(x_\alpha) f_\alpha(x_\alpha) dx_\alpha$ and define the centered B spline basis $b_{J,\alpha}(x_\alpha)$ and the standardized B spline basis $B_{J,\alpha}(x_\alpha)$ respectively as

$$b_{J,\alpha}(x_\alpha) = b_J(x_\alpha) - \frac{c_{J,\alpha}}{c_{J-1,\alpha}} b_{J-1}(x_\alpha),$$

$$B_{J,\alpha}(x_\alpha) = \frac{b_{J,\alpha}(x_\alpha)}{\left\{ \int b_{J,\alpha}^2(x_\alpha) f_\alpha(x_\alpha) dx_\alpha \right\}^{1/2}}, 1 \le J \le N + 1, \tag{A.4}$$

so that $\mathsf{E}\, B_{J,\alpha}(X_\alpha) \equiv 0$, $\mathsf{E}\, B_{J,\alpha}^2(X_\alpha) \equiv 1$.

One can rewrite with slight abuse of notations the log-likelihood $\widehat{L}(g)$ given in (7) as

$$\widehat{L}(g) = \widehat{L}(\boldsymbol{\lambda}) = n^{-1} \sum_{i=1}^n \left[ Y_i \boldsymbol{\lambda}^\mathsf{T} \mathbf{B}(\mathbf{X}_i) - b\left\{ \boldsymbol{\lambda}^\mathsf{T} \mathbf{B}(\mathbf{X}_i) \right\} \right], \tag{A.5}$$

with $g(\mathbf{X}_i) = \boldsymbol{\lambda}^\mathsf{T} \mathbf{B}(\mathbf{X}_i) \in G_n^0$, $\boldsymbol{\lambda} = (\lambda_0, \lambda_{J,\alpha})_{1 \le J \le N+1, 1 \le \alpha \le d}^\mathsf{T} \in \mathbb{R}^{N_d}$ with $N_d = (N+1)d + 1$, $\mathbf{B}(\mathbf{x}) = \{1, B_{1,1}(x_1), ..., B_{N+1,d}(x_d)\}^\mathsf{T}$ and $B_{J,\alpha}(x_\alpha)$ as given in (A.4). It is straightforward to verify that the gradient and Hessian of $\widehat{L}(\boldsymbol{\lambda})$ are

$$\nabla \widehat{L}(\boldsymbol{\lambda}) = n^{-1} \sum_{i=1}^n \left[ Y_i \mathbf{B}(\mathbf{X}_i) - b'\left\{ \boldsymbol{\lambda}^\mathsf{T} \mathbf{B}(\mathbf{X}_i) \right\} \mathbf{B}(\mathbf{X}_i) \right], \tag{A.6}$$

$$\nabla^2 \widehat{L}(\boldsymbol{\lambda}) = -n^{-1} \sum_{i=1}^n b''\left\{ \boldsymbol{\lambda}^\mathsf{T} \mathbf{B}(\mathbf{X}_i) \right\} \mathbf{B}(\mathbf{X}_i) \mathbf{B}(\mathbf{X}_i)^\mathsf{T}.$$

PROPOSITION A.1 *Under Assumptions (A1)-(A5) and (A7), for $m \in M$ with $M$ given in (25) and $\widehat{m}$ as in (8), as $n \to \infty$, $\|m - \widehat{m}\|_{2,n} + \|m - \widehat{m}\|_2 = \mathcal{O}_{a.s.}\left(N^{1/2} n^{-1/2} \log n\right)$ and $\|m - \widehat{m}\|_\infty = \mathcal{O}_{a.s.}\left(N n^{-1/2} \log n\right)$. With probability approaching 1, the Hessian matrix $\nabla^2 \widehat{L}(\boldsymbol{\lambda})$ satisfies that $\nabla^2 \widehat{L}(\boldsymbol{\lambda}) < \mathbf{0}, \forall \boldsymbol{\lambda}$ and $\nabla^2 \widehat{L}(\boldsymbol{\lambda}) \le -c_b c_V \mathbf{I}$ if $\boldsymbol{\lambda}^\mathsf{T} \mathbf{B}(\mathbf{X}_i) \in \Theta, 1 \le i \le n$.*

PROOF. See Lemma A.13 of LYH, Assumption (A2), equation (A.6) and Lemma A.11 of LYH.

## A.2. Proof of Theorem 1

PROOF. Define a stochastic process $\widehat{\varepsilon}_n(x_1) = n^{-1} \sum_{i=1}^n K_h(X_{i1} - x_1)\, \xi_i, x_1 \in [0,1]$ with $\xi_i$ given in (A.1), then (12) and (13) show that

$$\sup_{x_1 \in [h, 1-h]} \left| \widetilde{m}_{\mathrm{K},1}(x_1) - m_1(x_1) - D_1^{-1}(x_1)\widehat{\varepsilon}_n(x_1) \right| = \mathcal{O}_{a.s.}\left( h^2 + n^{-1/2} h^{1/2} \log n \right),$$

which, together with (11), lead to

$$\sup_{x_1 \in [h, 1-h]} \left| \widehat{m}_{\mathrm{SBK},1}(x_1) - m_1(x_1) - D_1^{-1}(x_1)\widehat{\varepsilon}_n(x_1) \right| \tag{A.7}$$
$$= \mathcal{O}_{a.s.}\left( h^2 + n^{-1/2} h^{1/2} \log n + n^{-1/2} \log n \right) = \mathcal{O}_{a.s.}\left( h^2 + n^{-1/2} \log n \right).$$

Using $v_1(x_1)$ given in (15), one can standardize $\widehat{\varepsilon}_n(x_1)$ and then replace $x_1$ by $t$ to obtain

$$\widehat{\zeta}_n(t) = (nh)^{1/2} v_1^{-1}(t) \widehat{\varepsilon}_n(t) = (nh)^{1/2} v_1^{-1}(t) \left\{ n^{-1} \sum_{i=1}^n K_h(X_{i1} - t)\, \xi_i \right\}. \tag{A.8}$$

Assumptions (A5), (A8) imply that the following Rosenblatt transformation to the 2-dimensional sequence $\{X_{1i}, \xi_i\}_{i=1}^n$ produces $\{X_{i1}', \xi_i'\}_{i=1}^n$ with $(X_{i1}', \xi_i')$ uniformly distributed on $[0,1]^2$:

$$(X_{i1}', \xi_i') = T(X_{1i}, \xi_i) = \left\{ F_{X_1}(X_{1i}), F_{\xi|X_1}(\xi_i|X_{1i}) \right\}.$$

Denote $Z_n(x_1, \xi) = \sqrt{n}\left\{ F_n(x_1, \xi) - F(x_1, \xi) \right\}$ where $F_n(x_1, \xi)$ is the empirical distribution of $\{X_{i1}, \xi_i\}_{i=1}^n$, one can rewrite $\widehat{\xi}_n(t)$ as

$$\widehat{\xi}_n(t) = h^{-1/2} v_1^{-1}(t) \int \int K\left\{ (x_1 - t)/h \right\} \xi dZ_n(x_1, \xi).$$

By the strong approximation theorem in Tusnady (1977), there exists a version of the 2-dimensional Brownian Bridge $B_n(x_1', \xi')$ such that

$$\sup_{x_1, s} \left| Z_n(x_1, \xi) - B_n\{ T(x_1, \xi) \} \right| = \mathcal{O}_{a.s.}\left( n^{-1/2} \log^2 n \right). \tag{A.9}$$

19

Applying standard techniques used in Bickel and Rosenblatt (1973), Härdle (1989), one can show that

$$\sup_{t\in[0,1]} \left|\widehat{\xi}_n\left(t\right) - M_h\left(t\right)/\|K\|_2^2\right| = \mathcal{O}_p\left\{\left(\log n\right)^{-1/2}\right\}, \tag{A.10}$$

for a version of the $M_h\left(t\right)$ given in (A.2). Similar result can be found in Xia (1998).

Furthermore, (A.7) and (A.8) imply that

$$\sup_{x_1\in[h,1-h]} \left|\sigma_n^{-1}\left(x_1\right)\left\{\widehat{m}_{\mathrm{SBK},1}\left(x_1\right) - m_1\left(x_1\right)\right\} - \widehat{\xi}_n\left(x_1\right)\right| = \mathcal{O}_{a.s.}\left(n^{1/2}h^{5/2} + h^{1/2}\log n\right), \tag{A.11}$$

with $\sigma_n\left(x\right)$ given in (19). Under Assumption (A6), which entails that $\left(-2\log h\right)^{1/2}$ is of the same order as $\left(\log n\right)^{1/2}$, (A.10) and (A.11) can show that

$$\sup_{x_1\in[h,1-h]}\left(-2\log h\right)^{1/2}\left|\sigma_n^{-1}\left(x_1\right)\left|\widehat{m}_{\mathrm{SBK},1}\left(x_1\right) - m_1\left(x_1\right)\right| - \left|M_h\left(t\right)\right|/\|K\|_2^2\right| \tag{A.12}$$

$$= \mathcal{O}_{a.s}\left\{\left(\log n\right)^{1/2} \times \left(n^{1/2}h^{5/2} + h^{1/2}\log n\right)\right\} + \mathcal{O}_p\left(1\right) = \mathcal{O}_p\left(1\right).$$

Finally, Theorem 1 follows from Lemma A.1 and Slutsky's Theorem.

## A.3. Proof of Theorem 2

Prior to proving Theorem 2, we restate Proposition A.1 for any index set $S \subset \{1, 2, \ldots d\}$.

Denote by $\boldsymbol{\lambda} = \left(\lambda_0, \lambda_{J,\alpha}\right)_{1\leq J\leq N+1, 1\leq\alpha\leq d}^{\mathsf{T}}$ an arbitrary vector. For any $S$, $N_S = 1 + (N+1)\,\#\left(S\right)$, denote

$$\boldsymbol{\lambda}_S = \left(\lambda_0, \lambda_{J,\alpha}\right)_{1\leq J\leq N+1, \alpha\in S}^{\mathsf{T}} \in \mathbb{R}^{N_S}, \mathbf{B}_S\left(\mathbf{x}\right) = \left\{1, B_{J,\alpha}\left(x_\alpha\right)\right\}_{1\leq J\leq N+1, \alpha\in S}^{\mathsf{T}}, \tag{A.13}$$

and with slight abuse of notations

$$\widehat{L}_S\left(\boldsymbol{\lambda}_S\right) = \widehat{L}_S\left\{\boldsymbol{\lambda}_S^{\mathsf{T}}\mathbf{B}_S\left(\mathbf{x}\right)\right\} = n^{-1}\sum_{i=1}^{n}\left[Y_i\boldsymbol{\lambda}_S^{\mathsf{T}}\mathbf{B}_S\left(\mathbf{X}_i\right) - b\left\{\boldsymbol{\lambda}_S^{\mathsf{T}}\mathbf{B}_S\left(\mathbf{X}_i\right)\right\}\right] \tag{A.14}$$

whose maximizer is $\widehat{m}_S = \widehat{\boldsymbol{\lambda}}_S^{\mathsf{T}}\mathbf{B}_S\left(\mathbf{x}\right)$.

PROPOSITION A.2 *Under Assumptions (A1)-(A5) and (A7), for $m_S \in M_S$ given in (28), $\widehat{m}_S$ in (A.14), as $n \to \infty$, $\|m_S - \widehat{m}_S\|_{2,n} + \|m_S - \widehat{m}_S\|_2 = \mathcal{O}_{a.s.}\left(N^{1/2}n^{-1/2}\log n\right)$ and $\|m_S - \widehat{m}_S\|_\infty = \mathcal{O}_{a.s.}\left(Nn^{-1/2}\log n\right)$.*

Next, we consider two cases "underfitting" and "overfitting" for the index set $S$ to establish Theorem 2.

**Definition:** if $S \supset S_0$ and $S \neq S_0$, then $S$ overfits, while $S$ is underfitting if $S_0 \cap S \neq S_0$ with $S_0$ given in Theorem 2. We shall show that $\lim_{n\to\infty} \mathrm{P}\left(\mathrm{BIC}_S - \mathrm{BIC}_{S_0} > 0\right) = 1$ in both situations.

PROOF. **I: overfitting, i.e.,** $S \supset S_0$ **and** $S \neq S_0$.

Let $\boldsymbol{\lambda}_{S_0 S} = \left\{ \boldsymbol{\lambda}_{S_0}, \left(\lambda_{J,\alpha'}\right)^{\mathsf{T}}_{1\leq J\leq N+1, \alpha'\in S\backslash S_0} \right\}$, $\widehat{\boldsymbol{\lambda}}_{S_0 S} = \left\{ \widehat{\boldsymbol{\lambda}}_{S_0}, \left(\lambda_{J,\alpha'}\right)^{\mathsf{T}}_{1\leq J\leq N+1, \alpha'\in S\backslash S_0} \right\}$ with $\lambda_{J,\alpha'} \equiv 0$ and $\widehat{\boldsymbol{\lambda}}_S$
( or $\widehat{\boldsymbol{\lambda}}_{S_0}$) as the MLE of (A.14) (or when $S = S_0$). Note that $\widehat{L}_S\left(\widehat{\boldsymbol{\lambda}}_{S_0 S}\right) = \widehat{L}_{S_0}\left(\widehat{\boldsymbol{\lambda}}_{S_0}\right)$.

Using Taylor's expansion, $\exists$ a vector $\widetilde{\boldsymbol{\lambda}}_S$ between $\widehat{\boldsymbol{\lambda}}_S$ and $\widehat{\boldsymbol{\lambda}}_{S_0 S}$, i.e., $\widetilde{\boldsymbol{\lambda}}_S = \mathbf{t}\widehat{\boldsymbol{\lambda}}_S + \left(\mathbf{I}_{Ns} - \mathbf{t}\right)\widehat{\boldsymbol{\lambda}}_{S_0 S}$ with a $N_s \times N_s$ diagonal matrix $\mathbf{t}$ whose diagonal elements are in $[0,1]$ s.t.

$$\widehat{L}_{S_0}\left(\widehat{\boldsymbol{\lambda}}_{S_0}\right) - \widehat{L}_S\left(\widehat{\boldsymbol{\lambda}}_S\right) = \widehat{L}_S\left(\widehat{\boldsymbol{\lambda}}_{S_0 S}\right) - \widehat{L}_S\left(\widehat{\boldsymbol{\lambda}}_S\right) \tag{A.15}$$
$$= \left(\widehat{\boldsymbol{\lambda}}_{S_0 S} - \widehat{\boldsymbol{\lambda}}_S\right)^{\mathsf{T}} \triangledown \widehat{L}_S\left(\widehat{\boldsymbol{\lambda}}_S\right) + \frac{1}{2}\left(\widehat{\boldsymbol{\lambda}}_{S_0 S} - \widehat{\boldsymbol{\lambda}}_S\right)^{\mathsf{T}} \triangledown^2 \widehat{L}_S\left(\widetilde{\boldsymbol{\lambda}}_S\right)\left(\widehat{\boldsymbol{\lambda}}_{S_0 S} - \widehat{\boldsymbol{\lambda}}_S\right).$$

Since $\triangledown\widehat{L}_S\left(\widehat{\boldsymbol{\lambda}}_S\right) = 0$ and $\triangledown^2\widehat{L}_S\left(\widetilde{\boldsymbol{\lambda}}_S\right)$ is given in (A.6), for $\widetilde{m}_S = \widetilde{\boldsymbol{\lambda}}_S^{\mathsf{T}}\mathbf{B}_S\left(\mathbf{x}\right)$, one has

$$\widehat{L}_{S_0}\left(\widehat{\boldsymbol{\lambda}}_{S_0}\right) - \widehat{L}_S\left(\widehat{\boldsymbol{\lambda}}_S\right) \tag{A.16}$$
$$= -\left(2n\right)^{-1}\sum_{i=1}^{n} b''\left\{\widetilde{\boldsymbol{\lambda}}_S^{\mathsf{T}}\mathbf{B}_S\left(\mathbf{X}_i\right)\right\}\left\{\left(\widehat{\boldsymbol{\lambda}}_{S_0 S} - \widehat{\boldsymbol{\lambda}}_S\right)^{\mathsf{T}}\mathbf{B}_S\left(\mathbf{X}_i\right)\right\}\left\{\left(\widehat{\boldsymbol{\lambda}}_{S_0 S} - \widehat{\boldsymbol{\lambda}}_S\right)^{\mathsf{T}}\mathbf{B}_S\left(\mathbf{X}_i\right)\right\}^{\mathsf{T}}$$
$$= -\left(2n\right)^{-1}\sum_{i=1}^{n} b''\left\{\widetilde{m}_S\left(\mathbf{X}_i\right)\right\}\left\{\widehat{m}_{S_0}\left(\mathbf{X}_i\right) - \widehat{m}_S\left(\mathbf{X}_i\right)\right\}^2.$$

By $m = m_{S_0} = m_S \in M_{S_0} \subset M_S$, Proposition A.2 implies that $\|\widehat{m}_{S_0} - m\|_{2,n} = \mathcal{O}_{a.s.}\left(N^{1/2}n^{-1/2}\log n\right)$ and $\|\widehat{m}_S - m\|_{2,n} = \mathcal{O}_{a.s.}\left(N^{1/2}n^{-1/2}\log n\right)$, thus

$$\|\widehat{m}_{S_0} - \widehat{m}_S\|_{2,n}^2 = \mathcal{O}_{a.s.}\left(Nn^{-1}\log^2 n\right). \tag{A.17}$$

Similarily, one has $\|\widetilde{m}_S - m\|_\infty = o_{a.s.}\left(1\right)$, which warrants for large $n$ that $\widetilde{m}_S \in \Theta$ with $\Theta$ given in Assumption (A.2), so (A.16) implies that

$$0 \geq \widehat{L}_{S_0}\left(\widehat{\boldsymbol{\lambda}}_{S_0}\right) - \widehat{L}_S\left(\widehat{\boldsymbol{\lambda}}_S\right) \geq -\frac{C_b}{2}\|\widehat{m}_{S_0} - \widehat{m}_S\|_{2,n}^2. \tag{A.18}$$

As a result, BIC given in (29) shows that

$$\mathrm{BIC}_S - \mathrm{BIC}_{S_0} = 2\left\{\widehat{L}_{S_0}\left(\widehat{\boldsymbol{\lambda}}_{S_0}\right) - \widehat{L}_S\left(\widehat{\boldsymbol{\lambda}}_S\right)\right\} + \frac{N_S - N_{S_0}}{n}\log^3 n \tag{A.19}$$
$$\geq -C_b\|\widehat{m}_{S_0} - \widehat{m}_S\|_{2,n}^2 + \left(N + 1\right)n^{-1}\log^3 n,$$

21

which implies by (A.17) that $\lim_{n\to\infty} P\left(\mathrm{BIC}_S - \mathrm{BIC}_{S_0} > 0\right) = 1$.

**II: underfitting, i.e., $S_0 \cap S \neq S_0$.**

Let $S' = S_0 \cup S$ and denote by $\widehat{\boldsymbol{\lambda}}_{S_0}$, $\widehat{\boldsymbol{\lambda}}_S$ and $\widehat{\boldsymbol{\lambda}}_{S'}$ the MLEs in (A.14) for $S_0, S$ and $S'$, respectively. Since $S'$ overfits $S_0$, similarly to (A.18), one has

$$0 \geq \widehat{L}_{S_0}\left(\widehat{\boldsymbol{\lambda}}_{S_0}\right) - \widehat{L}_{S'}\left(\widehat{\boldsymbol{\lambda}}_{S'}\right) \geq -\frac{C_b}{2}\left\|\widehat{m}_{S_0} - \widehat{m}_{S'}\right\|_{2,n}^2, \tag{A.20}$$

Define a set $\Theta_{S'} = \left\{\boldsymbol{\lambda}_{S'} : \boldsymbol{\lambda}_{S'}^{\mathsf{T}} \mathbf{B}_{S'}\left(\mathbf{X}_i\right) \in \Theta, 1 \leq i \leq n\right\}$. which is compact and convex in $\mathbb{R}^{N_S}$. By definition,

$$\max_{1 \leq i \leq n}\left|\widehat{\boldsymbol{\lambda}}_{S'}^{\mathsf{T}} \mathbf{B}_{S'}\left(\mathbf{X}_i\right) - m\left(\mathbf{X}_i\right)\right| \leq \left\|\widehat{m}_{S'} - m\right\|_{\infty} = \mathcal{O}_{a.s.}\left(Nn^{-1/2}\log n\right),$$

so for large $n$, with probability approaching 1, $\widehat{\boldsymbol{\lambda}}_{S'} \in \Theta_{S'}$, so Proposition A.1 ensures that, with probability approaching 1, the Hessian matrix $\nabla^2 \widehat{L}_{S'}\left(\widehat{\boldsymbol{\lambda}}_{S'}\right) \leq -c_b c_V \mathbf{I}_{N_{S'}}$, while $\nabla \widehat{L}_{S'}\left(\widehat{\boldsymbol{\lambda}}_{S'}\right) = 0$ and $\nabla^2 \widehat{L}_{S'}\left(\boldsymbol{\lambda}_{S'}\right) \leq \mathbf{0}$, $\forall \boldsymbol{\lambda}_{S'}$. Thus, with probability approaching 1, there exists a constant $c_1 > 0$ such that

$$\widehat{L}_{S'}\left(\boldsymbol{\lambda}_{S'}\right) - \widehat{L}_{S'}\left(\widehat{\boldsymbol{\lambda}}_{S'}\right) \leq \begin{cases} -2^{-1}c_b c_V \left\|\boldsymbol{\lambda}_{S'} - \widehat{\boldsymbol{\lambda}}_{S'}\right\|^2, & \text{if } \boldsymbol{\lambda}_{S'} \in \Theta_{S'} \\ \max\limits_{\boldsymbol{\lambda}_{S'} \in \partial\Theta_{S'}} \widehat{L}_S\left(\boldsymbol{\lambda}_S\right) - \widehat{L}_{S'}\left(\widehat{\boldsymbol{\lambda}}_{S'}\right) \leq -c_1, & \text{otherwise} \end{cases} . \tag{A.21}$$

Next, define a new vector $\widehat{\boldsymbol{\lambda}}_{SS'} = \left\{\widehat{\boldsymbol{\lambda}}_S, \left(\lambda_{J,\alpha'}\right)_{1 \leq J \leq N+1, \alpha' \in S' \setminus S}^{\mathsf{T}}\right\}$ with $\lambda_{J,\alpha'} \equiv 0$ and note that $\widehat{\boldsymbol{\lambda}}_{SS'}^{\mathsf{T}}$ $\mathbf{B}_{S'}\left(\mathbf{x}\right) \equiv \widehat{m}_S\left(\mathbf{x}\right), \widehat{\boldsymbol{\lambda}}_{S'}^{\mathsf{T}} \mathbf{B}_{S'}\left(\mathbf{x}\right) \equiv \widehat{m}_{S'}\left(\mathbf{x}\right)$, so applying Lemma A.5 of Wang and Yang (2007), there exists a constant $C_0 > 0$ such that

$$\left\|\widehat{\boldsymbol{\lambda}}_{SS'} - \widehat{\boldsymbol{\lambda}}_{S'}\right\|^2 \geq C_0^{-1}\left\|\widehat{m}_S - \widehat{m}_{S'}\right\|_2^2.$$

Applying Proposition A.2 entails that

$$\left|\left\|\widehat{m}_S - \widehat{m}_{S'}\right\|_2^2 - \left\|m_S - m_{S'}\right\|_2^2\right| = \mathcal{O}_{a.s.}\left(Nn^{-1}\overset{2}{\log} n\right)$$

while the definitions of underfitting and overfitting lead to

$$\left\|m_S - m_{S'}\right\|_2^2 = \left\|m_S - m\right\|_2^2 = c_S > 0$$

and thus

$$\left\|\widehat{\boldsymbol{\lambda}}_{SS'} - \widehat{\boldsymbol{\lambda}}_{S'}\right\|^2 \geq C_0^{-1}c_S + \mathcal{O}_{a.s.}\left(Nn^{-1}\log^2 n\right). \tag{A.22}$$

Note next that

$$\widehat{L}_S\left(\widehat{\boldsymbol{\lambda}}_S\right) - \widehat{L}_{S'}\left(\widehat{\boldsymbol{\lambda}}_{S'}\right) = \widehat{L}_{S'}\left(\widehat{\boldsymbol{\lambda}}_{SS'}\right) - \widehat{L}_{S'}\left(\widehat{\boldsymbol{\lambda}}_{S'}\right) \tag{A.23}$$

which according to (A.21), is bounded by

$$\leq \begin{cases} -2^{-1}c_b c_V \left\|\widehat{\boldsymbol{\lambda}}_{SS'} - \widehat{\boldsymbol{\lambda}}_{S'}\right\|^2, & \text{if } \widehat{\boldsymbol{\lambda}}_{SS'} \in \Theta_{S'} \\ \max_{\boldsymbol{\lambda}_{S'} \in \partial \Theta_{S'}} \widehat{L}_S\left(\boldsymbol{\lambda}_S\right) - \widehat{L}_{S'}\left(\widehat{\boldsymbol{\lambda}}_{S'}\right) \leq -c_1, & \text{otherwise} \end{cases}$$

which is, according to (A.22), bounded by

$$\begin{aligned} &\leq \max\left(-2^{-1}c_b c_V C_0^{-1} c_S, -c_1\right) + \mathcal{O}_{a.s.}\left(N n^{-1}\log^2 n\right) \\ &= -c_2 + \mathcal{O}_{a.s.}\left(N n^{-1}\log^2 n\right), \text{ for a constant } c_2 > 0. \end{aligned}$$

The above bound, together with (A.17), (A.20) and (A.23) lead to $\widehat{L}_{S_0}\left(\widehat{\boldsymbol{\lambda}}_{S_0}\right) - \widehat{L}_S\left(\widehat{\boldsymbol{\lambda}}_S\right)$

$$\begin{aligned} &= \left\{\widehat{L}_{S_0}\left(\widehat{\boldsymbol{\lambda}}_{S_0}\right) - \widehat{L}_S\left(\widehat{\boldsymbol{\lambda}}_{S'}\right)\right\} - \left\{\widehat{L}_S\left(\widehat{\boldsymbol{\lambda}}_S\right) - \widehat{L}_{S'}\left(\widehat{\boldsymbol{\lambda}}_{S'}\right)\right\} \\ &\geq c_2 + \mathcal{O}_{a.s.}\left(N n^{-1}\log^2 n\right). \end{aligned} \tag{A.24}$$

Finally, (A.24) implies that

$$\text{BIC}_S - \text{BIC}_{S_0} = 2\left\{\widehat{L}_{S_0}\left(\widehat{\boldsymbol{\lambda}}_{S_0}\right) - \widehat{L}_S\left(\widehat{\boldsymbol{\lambda}}_S\right)\right\} + \frac{N_S - N_{S_0}}{n}\log^3 n \geq c_2 + \mathcal{O}_p(1), \tag{A.25}$$

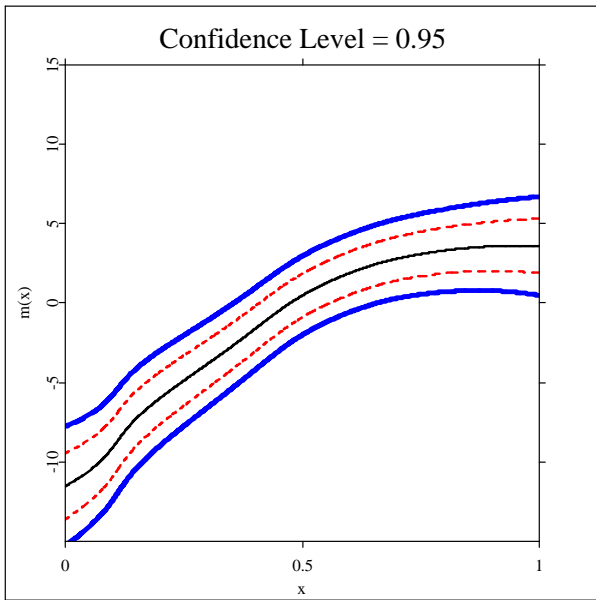and thus $\lim_{n \to \infty} \text{P}\left(\text{BIC}_S - \text{BIC}_{S_0} > 0\right) = 1$.

# References

[1] Berg, D. (2007), "Bankruptcy Prediction by Generalized Additive Models," *Applied Stochastic Models in Business and Industry*, 23, 129–143.

[2] Bernhardsen, E. (2001), *A Model of Bankruptcy Prediction*, Norges Bank WP, 2001.

[3] Bickel, P.J. and Rosenblatt, M. (1973), "On Some Global Measures of the Deviations of Density Function Estimates," *Annals of Statistics*, 1, 1071–1095.

[4] Engelmann, B., Hayden, E. and D. Tasche. (2003), "Testing Rating Accuracy," *Risk*, 16, 82–86.

[5] Fan, J., Härdle, W. and Mammen, E. (1998). Direct estimation of low-dimensional components in additive models. *Ann. Statist.* 26, 943–971.

[6] Fan, J. and Zhang, W. Y. (2000), "Simultaneous Confidence Bands and Hypothesis Testing in Varying-coefficient Models," *Scandinavian Journal of Statistics*, 27, 715–731.

[7] Härdle, W. (1989). Asymptotic maximal deviation of M-smoothers. *Journal of Multivariate Analysis* 29, 163–179.

[8] Härdle, W., Hoffmann, L. and Moro, R. (2011), *Learning Machines Supporting Bankruptcy prediction. Statistical Tools in Finance and Insurance* (2nd ed.), Cizek, Härdle, Weron, Springer Verlag.

[9] Hastie, T. J. and Tibshirani, R. J. (1990), *Generalized Additive Models*, Chapman and Hall, London.

[10] He, X., Fung, W. and Zhu, Z. (2005), "Robust Estimation in Generalized Partial Linear Models for Clustered Data," *Journal of the American Statistical Association*, 100, 1176–1184.

[11] He, X., Zhu, Z and Fung, W. (2002), "Estimation in A Semiparamtric Model for Longitudinal Data with Unspecified Dependence Structure," *Biometrika*, 89, 579–590.

[12] Horowitz, J. and Mammen, E. (2004), "Nonparametric Estimation of An Additive Model with A Link Function," *Annals of Statistics*, 32, 2412–2443.

[13] Horowitz, J., Klemelä, J. and Mammen, E. (2006), "Optimal Estimation in Additive Regression," *Bernoulli*, 12, 271–298.

[14] Huang, J. Z. and Yang, L. (2004), "Identification of Nonlinear Additive Autoregression Models," *Journal of the Royal Statistical Society: Series B*, 66, 463–477.

[15] Linton, O. B. (1997), "Efficient Estimation of Additive Nonparametric Regression Models," *Biometrika*, 84, 469–473.

[16] Linton, O. B. and Härdle, W. (1996), "Estimation of Additive Regression Models with Known Links," *Biometrika*, 83, 529–540.

[17] Liu, R. and Yang, L. (2010), "Spline-backfitted Kernel Smoothing of Additive Coefficient Model," *Econometric Theory*, 26, 29–59.
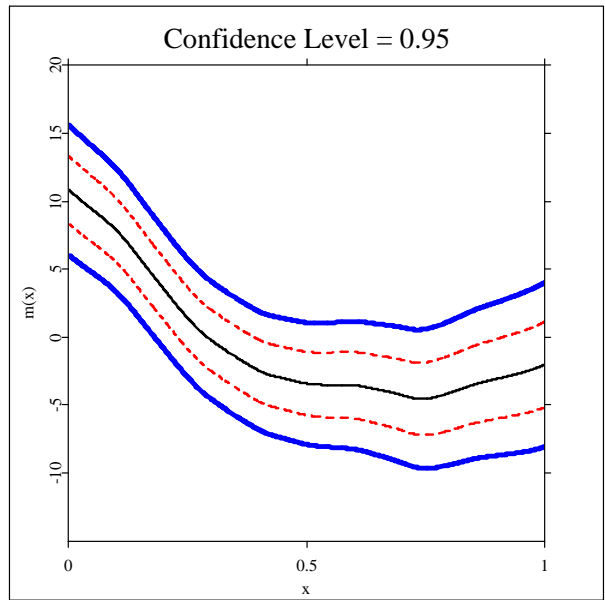
[18] Liu, R., Yang, L and Härdle, W. (2013), "Oracally Efficient Two-Step Estimation of Generalized Additive Model," *Journal of the American Statistical Association*, 108, 619–631.

[19] Ma, S. and Yang, L. (2011), "Spline-backfitted Kernel Smoothing of Partially Linear Additive model," *Journal of Statistical Planning and Inference*, 141, 204–219.

[20] Ma, S., Yang, L. and Carroll, R. J. (2012), " Simultaneous Confidence Band for Sparse Longitudinal Regression," *Statistica Sinica*, 22, 95–122.

[21] Ryser, M. and Denzler, S. (2009), "Selecting Credit Rating Models: A Cross-validation-based Comparison of Discriminatory Power," *Financ Mark Portf Manag*, 23, 187–203.

[22] Shina, Y. and Moore, W. (2003), "Explaining Credit Rating Differences between Japanese and U.S. Agencies," *Review of Financial Economics*, 12, 327–344.

[23] Song, Q. and Yang, L. (2010), "Oracally Efficient Spline Smoothing of Nonlinear Additive Autoregression Model with Simultaneous Confidence Band," *Journal of Multivariate Analysis*, 101, 2008–2025.

[24] Stone, C. J. (1985), "Additive Regression and Other Nonparametric Models," *Annals of Statistics*, 13, 689–705.

[25] Stone, C. J. (1986), "The Dimensionality Reduction Principle for Generalized Additive Models," *Annals of Statistics*, 14, 590–606

[26] Tusnady, G. (1977), "A Remark on the Approximation of the Sample Distribution Function in the Multidimensional Case," *Period. Math. Hungar.*, 8, 53–55.

[27] Wang, L. and Yang, L. (2007), "Spline-backfitted Kernel Smoothing of Nonlinear Additive Autoregression Model," *Annals of Statistics*, 35, 2474–2503.

[28] Wang, L., Li, H., and Huang, J. (2008), "Variable Selection in Nonparametric Varying-coefficient Models for Analysis of Repeated Measurements," *Journal of the American Statistical Association*, 103, 1556–1569.

[29] Wu, W. and Zhao, Z. (2007), "Inference of Trends in Time Series," *Journal of the Royal Statistical Society: Series B*, 69, 391–410.

[30] Xia, Y. (1998), "Bias-corrected Confidence Bands in Nonparametric Regression," *Journal of the Royal Statistical Society: Series B*, 60, 797–811.

[31] Xue, L. and Liang, H. (2010), "Polynomial Spline Estimation for A Generalized Additive Coefficient Model," *Scandinavian Journal of Statistics*, 37, 26–46.

[32] Xue, L. and Yang, L. (2006), "Additive Coefficient Modeling via Polynomial Spline," *Statistica Sinica*, 16, 1423–1446.

[33] Yang, L., Sperlich, S. and Härdle, W. (2003), "Derivative Estimation and Testing in Generalized Additive Models," *Journal of Statistical Planning and Inference*, 115, 521–542.

[34] Zhang, H. and Lin, Y. (2006), "Component Selection and Smoothing for Nonparametric Regression in Exponential Families," *Statistica Sinica*, 16, 1021–1042.

[35] Zhao, Z. and Wu, W. (2008), "Confidence Bands in Nonparametric Time Series Regression," *Annals of Statistics*, 36, 1854–1878.
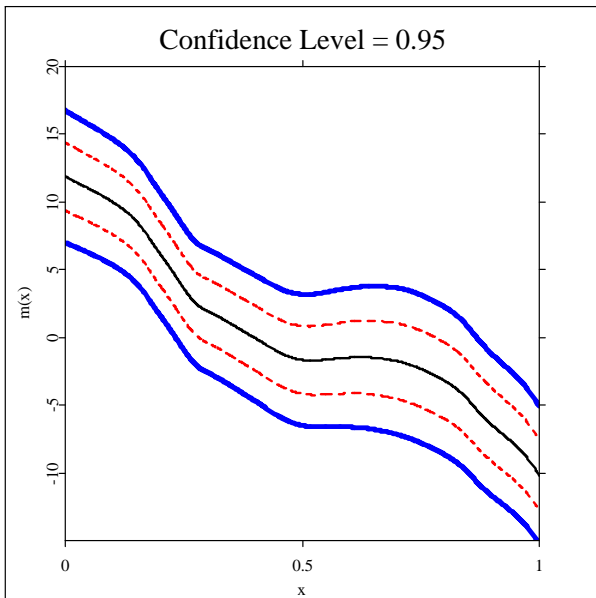
Figure 1: Plots of the rating factors in (a)-(c): SBK estimators (thin), 95% CIs (dashed) and 95% SCCs (thick). Plot of the CAPs in (d): Perfect (dashed), GAM (thick solid), GLM(thin solid), noninformative(dotted).
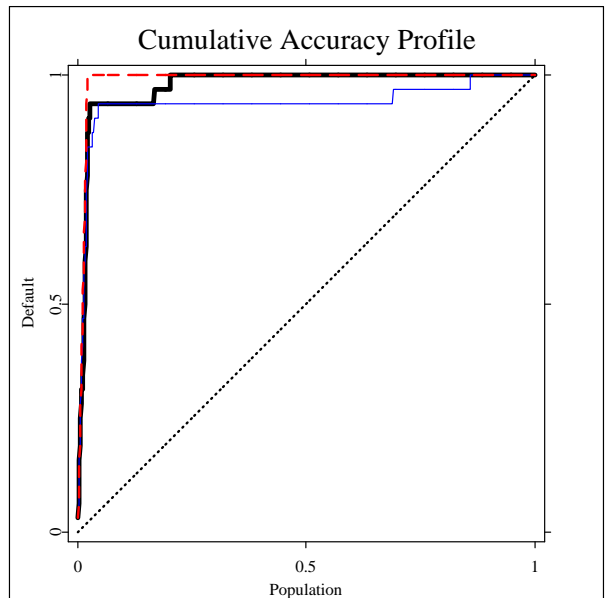
(a) **Current Liability**

(b) 3 **Months Earlier Return**

(c) 6 **Months Earlier Return**

(d) **The CAP Curves**

| d | r | n | Computing Time | | | Accuracy | | | | | |
| | | | BIC | COSSO | Ratio | BIC | | | COSSO | | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | | 250 | 0.43 | 5.90 | 13.7 | 4 | 90 | 6 | 20 | 66 | 14 |
| | 0 | 500 | 0.58 | 11.50 | 19.8 | 0 | 95 | 5 | 7 | 85 | 8 |
| | | 1000 | 1.05 | 27.89 | 26.6 | 0 | 99 | 1 | 6 | 90 | 4 |
| 10 | | | | | | | | | | | |
| | | 250 | 0.51 | 6.12 | 12.0 | 31 | 62 | 7 | 42 | 44 | 14 |
| | 0.5 | 500 | 0.64 | 12.90 | 20.1 | 3 | 91 | 6 | 17 | 73 | 10 |
| | | 1000 | 1.12 | 29.43 | 26.3 | 0 | 99 | 1 | 12 | 82 | 6 |
| | | 250 | 1.89 | — | — | 87 | 10 | 3 | — | — | — |
| | 0 | 500 | 2.76 | 209.34 | 76 | 14 | 77 | 9 | 21 | 46 | 33 |
| | | 1000 | 5.14 | 531.96 | 103 | 0 | 96 | 4 | 1 | 89 | 10 |
| 50 | | | | | | | | | | | |
| | | 250 | 2.02 | — | — | 90 | 4 | 6 | — | — | — |
| | 0.5 | 500 | 2.87 | 215.64 | 75 | 63 | 32 | 5 | 44 | 26 | 30 |
| | | 1000 | 5.39 | 545.43 | 101 | 10 | 86 | 4 | 10 | 71 | 19 |

Table 1: Simulation results for the proposed BIC method and COSSO with $d = 10$ and $50$. For each setup, the first, second, and third columns under Accuracy give respectively the frequencies of under fitting, correct fitting, and over fitting over 100 replications. Ratio is the computing time of COSSO over that of BIC. For $d = 50$ and $n = 250$, COSSO becomes unstable to the point of crashing.

| $r$ | $n$ | | | $\alpha$ | | |
|---|---|---|---|---|---|---|
| | | 3 | 4 | 5 | 6 | 7 |
| 0.0 | 250 | 0.925 | 0.920 | 0.920 | 0.930 | 0.925 |
| | 500 | 0.955 | 0.945 | 0.945 | 0.950 | 0.955 |
| | 1000 | 0.950 | 0.945 | 0.950 | 0.945 | 0.950 |
| 0.5 | 250 | 0.910 | 0.910 | 0.915 | 0.920 | 0.920 |
| | 500 | 0.935 | 0.925 | 0.930 | 0.930 | 0.935 |
| | 1000 | 0.945 | 0.945 | 0.940 | 0.950 | 0.945 |

Table 2: The 95% SCC global coverage for $m_a(x)$ , $\alpha = 3, ..., 7$ from 200 replications

# SFB 649 Discussion Paper Series 2014

For a complete list of Discussion Papers published by the SFB 649,
please visit http://sfb649.wiwi.hu-berlin.de.

001    "Principal Component Analysis in an Asymmetric Norm" by Ngoc Mai
       Tran, Maria Osipenko and Wolfgang Karl Härdle, January 2014.

002    "A Simultaneous Confidence Corridor for Varying Coefficient Regression
       with Sparse Functional Data" by Lijie Gu, Li Wang, Wolfgang Karl Härdle
       and Lijian Yang, January 2014.

003    "An Extended Single Index Model with Missing Response at Random" by
       Qihua Wang, Tao Zhang, Wolfgang Karl Härdle, January 2014.

004    "Structural Vector Autoregressive Analysis in a Data Rich Environment:
       A Survey" by Helmut Lütkepohl, January 2014.

005    "Functional stable limit theorems for efficient spectral covolatility
       estimators" by Randolf Altmeyer and Markus Bibinger, January 2014.

006    "A consistent two-factor model for pricing temperature derivatives" by
       Andreas Groll, Brenda López-Cabrera and Thilo Meyer-Brandis, January
       2014.

007    "Confidence Bands for Impulse Responses: Bonferroni versus Wald" by
       Helmut Lütkepohl, Anna Staszewska-Bystrova and Peter Winker, January
       2014.

008    "Simultaneous Confidence Corridors and Variable Selection for
       Generalized Additive Models" by Shuzhuan Zheng, Rong Liu, Lijian Yang
       and Wolfgang Karl Härdle, January 2014.