

Adaptive Varying-Coefficient Linear Models

Jianqing Fan*

Department of Statistics
University of North Carolina
Chapel Hill, NC 27599, USA

Qiwei Yao†

Department of Statistics
London School of Economics
Houghton Street, London WC2A 2AE, UK

Zongwu Cai‡

Department of Mathematics
University of North Carolina
Charlotte, NC 28223, USA

August 16, 2002

Abstract

Varying-coefficient linear models arise from multivariate nonparametric regression, nonlinear time series modeling and forecasting, functional data analysis, longitudinal data analysis, and others. It has been a common practice to assume that the varying coefficients are functions of a given variable which is often called an *index*. To enlarge the modeling capacity substantially, this paper explores a class of varying-coefficient linear models in which the index is unknown and is estimated as a linear combination of regressors and/or other variables. We search for the index such that the derived varying-coefficient model provides the least-squares approximation to the underlying unknown multi-dimensional regression function. The search is implemented through a newly proposed hybrid back-fitting algorithm. The core of the algorithm is the alternating iteration between estimating the index through a one-step scheme and estimating coefficient functions through a one-dimensional local linear smoothing. The locally significant variables are selected in terms of a combined use of *t*-statistic and the Akaike information criterion. We further extend the algorithm for the models with two indices. Simulation shows that the proposed methodology has appreciable flexibility to model complex multivariate nonlinear structure and is practically feasible with average modern computers. The methods are further illustrated through the Canadian mink-muskrat data in 1925-1994 and the pound/dollar exchange rates in 1974-1983.

Keywords: Akaike information criterion; back-fitting algorithm; generalized cross-validation; local linear regression; local significant variable selection; one-step estimation; smoothing index.

*Supported partially by NSF grant DMS-0196041

†Supported partially by EPSRC Grant L16385 and BBSRC/EPSRC Grant 96/MMI09785.

‡Supported partially by NSF grant DMS-0072400 and funds provided by the University of North Carolina at Charlotte.

1 Introduction

Suppose that we are interested in estimating multivariate regression function $G(\mathbf{x}) \equiv E(Y|\mathbf{X} = \mathbf{x})$, where Y is a random variable and \mathbf{X} is a $d \times 1$ random vector. In this paper, we propose to *approximate* the regression function $G(\mathbf{x})$ by a varying-coefficient model

$$g(\mathbf{x}) = \sum_{j=0}^d g_j(\boldsymbol{\beta}^T \mathbf{x}) x_j, \quad (1.1)$$

where $\boldsymbol{\beta} \in \mathbb{R}^d$ is an unknown direction, $\mathbf{x} = (x_1, \dots, x_d)^T$, $x_0 = 1$, and coefficients $g_0(\cdot), \dots, g_d(\cdot)$ are unknown functions. We choose the direction $\boldsymbol{\beta}$ and coefficient functions $g_j(\cdot)$ such that $E\{G(\mathbf{X}) - g(\mathbf{X})\}^2$ is minimised. The appeal of this model is that once $\boldsymbol{\beta}$ is given, we can directly estimate $g_j(\cdot)$ by the standard one-dimensional kernel regression localised around $\boldsymbol{\beta}^T \mathbf{x}$. Furthermore, the coefficient functions $g_j(\cdot)$ can be easily displayed graphically, which may be particularly helpful to visualise how the surface $g(\cdot)$ changes. Model (1.1) appears linear in each coordinate of \mathbf{x} when the index $\boldsymbol{\beta}^T \mathbf{x}$ is fixed. It may include quadratic and cross-product terms of x_j (or more generally any given functions of x_j) as ‘new’ components of \mathbf{x} . Hence it has considerable flexibility to cater to complex multivariate nonlinear structure.

We develop an efficient back-fitting algorithm to estimate $g(\cdot)$. The virtue of the algorithm is the alternating iteration between estimating $\boldsymbol{\beta}$ through a one-step estimation scheme (Bickel, 1975) and estimating functions $g_j(\cdot)$ through a one-dimensional local linear smoothing. Since we apply smoothing on a scalar $\boldsymbol{\beta}^T \mathbf{X}$ only, the method suffers little from the so-called ‘curse of dimensionality’ which is the innate difficulty associated with multivariate nonparametric fittings. The generalized cross-validation method (GCV) for bandwidth selection is incorporated into the algorithm in an efficient manner. To avoid over-fitting, we delete local insignificant variables in terms of a combined use of t -statistic and the Akaike information criterion (AIC) which is adopted for its computational efficiency. The deletion of insignificant variables is particularly important when we include, for example, quadratic functions of x_j as new components in the model, which could lead to over-parametrisation. The proposed method has been further extended to estimate varying-coefficient models with two indices, one of which is known.

Varying coefficient models arise from various statistical contexts in slightly different forms. The vast amount of literature includes, among others, Cleveland, Grosse and Shyu (1992), Hastie and Tibshirani (1993), Carroll, Ruppert and Welsh (1998), Kauermann and Tutz (1999), Xia and Li (1999a), Zhang and Lee (1999, 2000), and Fan and Zhang (1999, 2000b) on local multi-dimensional regression; Ramsay and Silverman (1997) on functional data analysis; Hoover *et al.* (1998), Wu, Chiang and Hoover (1998), and Fan and Zhang (2000a) on longitudinal data analysis; Nicholls and Quinn (1982), Chen and Tsay (1993), and Cai, Fan and Yao (2000) on nonlinear time series; and Cai, Fan and Li (2000) on generalized linear models with varying coefficients. The form of model (1.1) is not new. It was proposed in Ichimura (1993). Recently, Xia and Li (1999b) extended the idea and the results of Härdle, Hall and Ichimura (1993) from the single-index model to the

adaptive varying-coefficient model (1.1). They proposed to estimate the coefficient functions with a given bandwidth and a direction β , and then to choose the bandwidth and the direction by cross-validation. Some theoretical results were derived under the assumption that the bandwidth was of the order $O(n^{-1/5})$ and the direction β was within an $O_p(n^{-1/2})$ consistent neighbourhood of the true value. However, the approach suffers from heavy computational expense. This somehow explains why most previous work assumed a known direction β . The new approach in this paper differs from those in three key aspects: (a) only a one-dimensional smoother is used in estimation, (b) the index coefficient β is estimated by data and (c) within a local region around $\beta^T \mathbf{x}$, we select significant variables x_j to avoid overfitting. Aspect (b) is different from Härdle, Hall and Ichimura (1993) and Xia and Li (1999b) since we estimate the coefficient functions and the direction simultaneously; no cross-validation is needed. This idea is similar *in spirit* to that of Carroll *et al.* (1997) who showed that a semiparametric efficient estimator of the direction β can be obtained. Further we provide a theorem (i.e. Theorem 1(ii) in section 2 below) on the model identification problem of the form (1.1), which has not been addressed before.

The rest of the paper is organised as follows. Section 2 deals with the adaptive varying-coefficient model (1.1). The extension to the case with two indices is outlined on section 3. The numerical results of two simulated examples are reported in section 4.1, which demonstrate that the proposed methodology is able to capture complex nonlinear structure with moderate sample sizes, and further the required computation typically takes less than a minute on a Pentium II 350MHz PC. The methodology is further illustrated in section 4.2 through Canadian mink-muskrat data in 1925-1944 and the pound/dollar exchange rates in 1974-1983. The technical proofs are relegated to the Appendix.

2 Adaptive varying-coefficient linear models

2.1 Approximation and identifiability

Since $G(\mathbf{x}) = E(Y | \mathbf{X} = \mathbf{x})$ is a conditional expectation, it holds that

$$E\{Y - g(\mathbf{X})\}^2 = E\{Y - G(\mathbf{X})\}^2 + E\{G(\mathbf{X}) - g(\mathbf{X})\}^2$$

for any $g(\cdot)$. Therefore, the search for the least-squares approximation $g(\cdot)$ of $G(\cdot)$, as defined in (1.1), is equivalent to the search for such a $g(\cdot)$ that $E\{Y - g(\mathbf{X})\}^2$ obtains the minimum. Theorem 1(i) below indicates that there always exists such a $g(\cdot)$ under mild conditions. Obviously, if $G(\mathbf{x})$ is in the form of the RHS of (1.1), $g(\mathbf{x}) \equiv G(\mathbf{x})$. The second part of the theorem points out that the coefficient vector β is unique up to a constant unless $g(\cdot)$ is in a class of special quadratic functions (see (2.2) below). In fact, model (1.1) is an over-parametrised form in the sense that one of functions $g_j(\cdot)$ can be represented in terms of the others. Theorem 1(ii) confirms that once the direction β is specified, the function $g(\cdot)$ has a representation with at most d (instead of $d + 1$) $g_j(\cdot)$. Furthermore, those $g_j(\cdot)$ are identifiable.

Theorem 1. (i) Assume that the distribution function of (\mathbf{X}, Y) is continuous, and $E\{Y^2 + \|\mathbf{X}\|^2\} < \infty$. Then, there exists a $g(\cdot)$ defined by (1.1) for which

$$E\{Y - g(\mathbf{X})\}^2 = \inf_{\boldsymbol{\alpha}} \inf_{f_0, \dots, f_d} E \left\{ Y - \sum_{j=0}^d f_j(\boldsymbol{\alpha}^T \mathbf{X}) X_j \right\}^2, \quad (2.1)$$

where the first infimum is taken over all unit vectors in \Re^d , and the second over all measurable functions $f_0(\cdot), \dots, f_d(\cdot)$.

(ii) For any given twice differentiable $g(\cdot)$ of the form (1.1), if we choose $\|\boldsymbol{\beta}\| = 1$, and the first non-zero component of $\boldsymbol{\beta}$ positive, such a $\boldsymbol{\beta}$ is unique unless $g(\cdot)$ is of the form

$$g(\mathbf{x}) = \boldsymbol{\alpha}^T \mathbf{x} \boldsymbol{\beta}^T \mathbf{x} + \boldsymbol{\gamma}^T \mathbf{x} + c, \quad (2.2)$$

where $\boldsymbol{\alpha}, \boldsymbol{\gamma} \in \Re^d$, $c \in \Re$ are constants, and $\boldsymbol{\alpha}$ and $\boldsymbol{\beta}$ are not in parallel with each other. Furthermore, once $\boldsymbol{\beta} = (\beta_1, \dots, \beta_d)^T$ is given and $\beta_d \neq 0$, we may let $g_d(\cdot) \equiv 0$. Consequently, all the other $g_j(\cdot)$ are uniquely determined.

Remark 1. If the conditional expectation $G(\mathbf{x}) = E(Y|\mathbf{X} = \mathbf{x})$ cannot be expressed in the form of the RHS of (1.1), there may exist more than one $g(\mathbf{x})$ of the form of (1.1), for which (2.1) holds. For example, let $Y = X_1^2 + X_2^2$, where both X_1 and X_2 are independent random variables uniformly distributed on $[0, 1]$. Then $G(x_1, x_2) = x_1^2 + x_2^2$, which is not in the form of (1.1). However, (2.1) holds for both $g(x_1, x_2) = 1.25x_1^2$, and $1.25x_2^2$.

Without loss of the generality, we always assume from now on that in model (1.1), $\|\boldsymbol{\beta}\| = 1$ and the first non-zero component of $\boldsymbol{\beta}$ is positive. To avoid the complication caused by the lack of uniqueness of the index direction $\boldsymbol{\beta}$, we always assume that $G(\cdot)$ admits a unique least-squares approximation of $g(\cdot)$ which cannot be expressed in the form of (2.2).

2.2 Estimation

Let $\{(\mathbf{X}_t, Y_t); 1 \leq t \leq n\}$ be observations from a strictly stationary process with the same marginal distribution as (\mathbf{X}, Y) . In order to estimate the surface $g(\cdot)$ defined by (1.1) and (2.1), we need to search for the minimiser of $\{f_j(\cdot)\}$ for any given direction $\boldsymbol{\alpha}$ and then find the direction at which the mean squared error (MSE) is minimised. An exhaustive search is intractable. We develop a back-fitting algorithm for this optimisation problem.

Let $\beta_d \neq 0$. It follows from Theorem 1(ii) that we only search for an approximation in the form

$$g(\mathbf{x}) = \sum_{j=0}^{d-1} g_j(\boldsymbol{\beta}^T \mathbf{x}) x_j. \quad (2.3)$$

Our task can be formally split into two parts — estimation of functions $g_j(\cdot)$ with a given $\boldsymbol{\beta}$ and estimation of the index coefficient $\boldsymbol{\beta}$ with given $g_j(\cdot)$. We also discuss how to choose the smoothing parameter h in terms of GCV (Wahba, 1977), and how to apply backward deletion to choose locally significant variables in terms of a combined use of t -statistic and AIC. The algorithm for practical implementation will be summarised at the end of this section.

The computer-intensive nature of the problem prevents us from exploring more sophisticated methods which may lead to an improvement in performance at the cost of computing time. For example, various plug-in methods (Chapter 4 of Fan and Gijbels, 1996), a thorough GCV (see Step 3 in section 2.3 below) or the corrected AIC (Hurvich, Simonoff and Tsai, 1998; Simonoff and Tsai, 1999) would lead to better bandwidth selectors. The local variable selection could also be based solely on the AIC or the corrected AIC.

2.2.1 Local linear estimators for the $g_j(\cdot)$ with given β

For given β with $\beta_d \neq 0$, we need to estimate

$$g(\mathbf{X}) = \arg \min_{f \in \mathcal{F}(\beta)} E \left[\{Y - f(\mathbf{X})\}^2 \mid \beta^T \mathbf{X} \right], \quad (2.4)$$

where

$$\mathcal{F}(\beta) = \left\{ f(\mathbf{x}) = \sum_{j=0}^{d-1} f_j(\beta^T \mathbf{x}) x_j \mid f_0(\cdot), \dots, f_{d-1}(\cdot) \text{ measurable, and } E\{f(\mathbf{X})\}^2 < \infty \right\}.$$

The least-squares property in (2.4) leads to the estimators $\hat{g}_j(z) = \hat{b}_j$, $j = 0, \dots, d-1$, where $(\hat{b}_0, \dots, \hat{b}_{d-1})$ is the minimiser of the sum of weighted squares

$$\sum_{t=1}^n \left\{ Y_t - \sum_{j=0}^{d-1} b_j X_{tj} \right\}^2 K_h(\beta^T \mathbf{X}_t - z) w(\beta^T \mathbf{X}_t),$$

where $w(\cdot)$ is a bounded weight function with a bounded support, which is introduced to control the boundary effect. Note that only a one-dimensional kernel smoothing is used here.

The above estimation procedure is based on the local constant approximation: $g_j(y) \approx g_j(z)$ for y in a neighbourhood of z . Since local constant regression has several drawbacks compared with local linear regression (Fan and Gijbels, 1996), we consider the local linear estimators for the functions $g_0(\cdot), \dots, g_{d-1}(\cdot)$. This leads to minimising the sum

$$\sum_{t=1}^n \left[Y_t - \sum_{j=0}^{d-1} \{b_j + c_j(\beta^T \mathbf{X}_t - z)\} X_{tj} \right]^2 K_h(\beta^T \mathbf{X}_t - z) w(\beta^T \mathbf{X}_t) \quad (2.5)$$

with respect to $\{b_j\}$ and $\{c_j\}$. Define $\hat{g}_j(z) = \hat{b}_j$ and $\hat{g}_j(z) = \hat{c}_j$ for $j = 0, \dots, d-1$ and set

$$\hat{\boldsymbol{\theta}} \equiv (\hat{b}_0, \dots, \hat{b}_{d-1}, \hat{c}_0, \dots, \hat{c}_{d-1})^T.$$

It follows from the least-squares theory that

$$\hat{\boldsymbol{\theta}} = \boldsymbol{\Sigma}(z) \mathcal{X}^T(z) \mathcal{W}(z) \mathcal{Y}, \quad \text{and} \quad \boldsymbol{\Sigma}(z) = \{\mathcal{X}^T(z) \mathcal{W}(z) \mathcal{X}(z)\}^{-1}, \quad (2.6)$$

where $\mathcal{Y} = (Y_1, \dots, Y_n)^T$, $\mathcal{W}(z)$ is an $n \times n$ diagonal matrix with $K_h(\beta^T \mathbf{X}_i - z) w(\beta^T \mathbf{X}_i)$ as its i -th diagonal element, $\mathcal{X}(z)$ is an $n \times 2d$ matrix with $(\mathbf{U}_i^T, (\beta^T \mathbf{X}_i - z) \mathbf{U}_i^T)$ as its i -th row, and $\mathbf{U}_i = (1, X_{i1}, \dots, X_{i,d-1})^T$.

2.2.2 Search for β -direction with the $g_j(\cdot)$ fixed

The minimisation property in (2.1) suggests that we should search for β to minimise

$$R(\beta) = \frac{1}{n} \sum_{t=1}^n \left\{ Y_t - \sum_{j=0}^{d-1} g_j(\beta^T \mathbf{X}_t) X_{tj} \right\}^2 w(\beta^T \mathbf{X}_t), \quad (2.7)$$

We employ a one-step estimation scheme (see, for example, Bickel, 1975) to estimate β , which is in the spirit of the one-step Newton-Raphson estimation. We anticipate that the derived estimator is good if the initial value is reasonable (see Fan and Chen, 1999).

Suppose that $\hat{\beta}$ is the minimiser of (2.7). Then $\dot{R}(\hat{\beta}) = 0$, where $\dot{R}(\cdot)$ denotes the derivative of $R(\cdot)$. For any $\beta^{(0)}$ close to $\hat{\beta}$, we have the approximation

$$0 = \dot{R}(\hat{\beta}) \approx \dot{R}(\beta^{(0)}) + \ddot{R}(\beta^{(0)}) (\hat{\beta} - \beta^{(0)}),$$

where $\ddot{R}(\cdot)$ is the Hessian matrix of $R(\cdot)$. This leads to the one-step iterative estimator

$$\beta^{(1)} = \beta^{(0)} - \left\{ \ddot{R}(\beta^{(0)}) \right\}^{-1} \dot{R}(\beta^{(0)}), \quad (2.8)$$

where $\beta^{(0)}$ is the initial value. We re-scale $\beta^{(1)}$ such that it has unit norm with first non-vanishing element positive. It is easy to see from (2.7) that

$$\begin{aligned} \dot{R}(\beta) &= -\frac{2}{n} \sum_{t=1}^n \left\{ Y_t - \sum_{j=0}^{d-1} g_j(\beta^T \mathbf{X}_t) X_{tj} \right\} \left\{ \sum_{j=0}^{d-1} \dot{g}_j(\beta^T \mathbf{X}_t) X_{tj} \right\} \mathbf{X}_t w(\beta^T \mathbf{X}_t), \\ \ddot{R}(\beta) &= \frac{2}{n} \sum_{t=1}^n \left\{ \sum_{j=0}^{d-1} \dot{g}_j(\beta^T \mathbf{X}_t) X_{tj} \right\}^2 \mathbf{X}_t \mathbf{X}_t^T w(\beta^T \mathbf{X}_t) \\ &\quad - \frac{2}{n} \sum_{t=1}^n \left\{ Y_t - \sum_{j=0}^{d-1} g_j(\beta^T \mathbf{X}_t) X_{tj} \right\} \left\{ \sum_{j=0}^{d-1} \ddot{g}_j(\beta^T \mathbf{X}_t) X_{tj} \right\} \mathbf{X}_t \mathbf{X}_t^T w(\beta^T \mathbf{X}_t). \end{aligned} \quad (2.9)$$

In the above derivation, the derivative of the weight function $w(\cdot)$ is assumed to be 0 for the sake of simplicity. In practice, we usually let $w(\cdot)$ be an indicator function.

In case that the matrix $\ddot{R}(\cdot)$ is singular or nearly so, we adopt a ridge regression (Seifert and Gasser, 1996) as follows: use the estimator (2.8) with \ddot{R} replaced by \ddot{R}_r which is defined by the RHS of (2.9) with $\mathbf{X}_t \mathbf{X}_t^T$ replaced by $\mathbf{X}_t \mathbf{X}_t^T + q_n \mathbf{I}_d$ for some positive ridge parameter q_n .

Now we briefly state two alternative methods for estimating β , although they may not as efficient as the above method. The first one is based on a random search method, which is more direct and tractable when d is small. The basic idea is to keep drawing β randomly from the d -dimensional unit sphere and then computing $R(\beta)$. Stop the algorithm if the minimum fails to decrease significantly in, say, every 100 new draws. The second approach is to adapt the average derivative method of Newey and Stoker (1993) and Samarov (1993). Under model (1.1), the direction β is parallel to the expected difference between gradient vector of the regression surface and $(g_1(\beta^T \mathbf{x}), \dots, g_{d-1}(\beta^T \mathbf{x}), 0)^T$ and hence can be estimated by the average derivative method via iteration.

2.2.3 Bandwidth selection

We apply the generalized cross-validation (GCV) method, proposed by Wahba (1977) and Craven and Wahba(1979), to choose the bandwidth h in estimation of $\{g_j(\cdot)\}$. The criterion can be described as follows. For given β , let $\hat{Y}_t = \sum_{j=0}^{d-1} \hat{g}_j(\beta^T \mathbf{X}_t) X_{tj}$. It is easy to see that all those predicted values are in fact the linear combinations of $\mathcal{Y} = (Y_1, \dots, Y_n)^T$ with coefficients depending on $\{\mathbf{X}_t\}$ only. Namely, we may write

$$(\hat{Y}_1, \dots, \hat{Y}_n)^T = \mathbf{H}(h)\mathcal{Y},$$

where $\mathbf{H}(h)$ is the $n \times n$ hat matrix, independent of \mathcal{Y} . The GCV method selects h minimising

$$\text{GCV}(h) \equiv \frac{1}{n\{1 - n^{-1}\text{tr}(\mathbf{H}(h))\}^2} \sum_{t=1}^n \{Y_t - \hat{Y}_t\}^2 w(\beta^T \mathbf{X}_t),$$

which in fact is an estimate of the weighted mean integrated square errors. Under some regularity conditions, it holds that

$$\text{GCV}(h) = a_0 + a_1 h^4 + \frac{a_2}{nh} + o_p(h^4 + n^{-1}h^{-1}).$$

Thus, up to first order asymptotics, the optimal bandwidth is $h_{opt} = \{a_2/(4na_1)\}^{1/5}$. The coefficients of a_0 and a_1 and a_2 will be estimated from $\text{GCV}(h_k)$ via least-squares regression. This bandwidth selection rule, inspired by the empirical bias method of Ruppert (1997), will be applied outside the loops between estimating β and $\{g_j(\cdot)\}$. See section 2.2.5.

To calculate $\text{tr}\{\mathbf{H}(h)\}$, we note that for $1 \leq i \leq n$,

$$\hat{Y}_i = \frac{1}{n} \sum_{t=1}^n Y_t K_h(\beta^T \mathbf{X}_t - \beta^T \mathbf{X}_i) w(\beta^T \mathbf{X}_t) (\mathbf{U}_t^T, \mathbf{0}^T) \Sigma(\beta^T \mathbf{X}_i) \begin{pmatrix} \mathbf{U}_t \\ \mathbf{U}_t \frac{\beta^T (\mathbf{X}_t - \mathbf{X}_i)}{h} \end{pmatrix},$$

where $\mathbf{0}$ denotes the $d \times 1$ vector with all components 0, and $\Sigma(\cdot)$ is defined as in (2.6). The coefficient of Y_i on the RHS of the above expression is

$$\gamma_i \equiv \frac{1}{n} K_h(0) w(\beta^T \mathbf{X}_i) (\mathbf{U}_i^T, \mathbf{0}^T) \Sigma(\beta^T \mathbf{X}_i) \begin{pmatrix} \mathbf{U}_i \\ \mathbf{0} \end{pmatrix}.$$

Now, we have that $\text{tr}\{\mathbf{H}(h)\} = \sum_{i=1}^n \gamma_i$.

2.2.4 Choosing locally significant variables

As we discussed before, model (2.3) can be over-parametrised. Thus, it is necessary to select significant variables for each given z after an initial fitting. In our implementation, we use a backward stepwise deletion technique which relies on a modified AIC and t -statistics. More precisely, we delete the least significant variable in a given model according to t -values, which yields a new and reduced model. We select the best model according to AIC.

We start with the full model

$$g(\mathbf{x}) = \sum_{j=0}^{d-1} g_j(\beta^T \mathbf{x}) x_j. \quad (2.10)$$

For fixed $\beta^T \mathbf{X} = z$, (2.10) could be viewed as a (local) linear regression model. The least-squares estimator $\hat{\boldsymbol{\theta}} \equiv \hat{\boldsymbol{\theta}}(z)$ given in (2.6) entails

$$\text{RSS}_d(z) = \sum_{t=1}^n \left[Y_t - \sum_{j=0}^{d-1} \{ \hat{g}_j(z) + \hat{g}_j(z)(\beta^T \mathbf{X}_t - z) \} X_{tj} \right]^2 K_h(\beta^T \mathbf{X}_t - z) w(\beta^T \mathbf{X}_t).$$

The ‘degrees of freedom’ of $\text{RSS}_d(z)$ is $m(d, z) = n_z - p(d, z)$ where $n_z = \text{tr}\{\mathcal{W}(z)\}$ may be viewed as the number of observations used in the local estimation and $p(d, z) = \text{tr}\{\boldsymbol{\Sigma}(z)\mathcal{X}^T(z)\mathcal{W}^2(z)\mathcal{X}(z)\}$ as the number of local parameters. Now we define the AIC for this model as

$$\text{AIC}_d(z) = \log\{\text{RSS}_d(z)/m(d, z)\} + 2p(d, z)/n_z.$$

To delete the least significant variable among x_0, x_1, \dots, x_{d-1} , we search for x_k such that both $g_k(z)$ and $\dot{g}_k(z)$ are close to 0. The t -statistics for those two variables in the (local) linear regression are

$$t_k(z) = \frac{\hat{g}_k(z)}{\sqrt{c_k(z)\text{RSS}(z)/m(d, z)}} \quad \text{and} \quad t_{d+k} = \frac{\hat{\dot{g}}_k(z)}{\sqrt{c_{d+k}(z)\text{RSS}(z)/m(d, z)}}$$

respectively, where $c_k(z)$ is the $(k+1, k+1)$ -th element of matrix $\boldsymbol{\Sigma}(z)\mathcal{X}^T(z)\mathcal{W}^2(z)\mathcal{X}(z)\boldsymbol{\Sigma}(z)$. Discarding a common factor, we define

$$T_k^2(z) = \{\hat{g}_k(z)\}^2/c_k(z) + \{\hat{\dot{g}}_k(z)\}^2/c_{d+k}(z).$$

Let j be the minimiser of $T_k^2(z)$ over $0 \leq k < d$. We delete x_j from the full model (2.10). This leads to a model with $(d-1)$ ‘linear terms’. Repeating the above process, we may define $\text{AIC}_l(z)$ for all $1 \leq l \leq d$. If $\text{AIC}_k(z) = \min_{1 \leq l \leq d} \text{AIC}_l(z)$, the selected model should have $k-1$ ‘linear terms’ x_j .

2.3 Implementation

Now we outline the algorithm as follows.

Step 1: Standardise the data $\{\mathbf{X}_t\}$ such that it has sample mean 0 and sample covariance matrix \mathbf{I}_d . Specify an initial value of β .

Step 2: For each prescribed bandwidth value h_k , $k = 1, \dots, q$, repeat (a) and (b) below until two successive values of $R(\beta)$ defined in (2.7) differ insignificantly:

(a) For a given β , estimate $g_j(\cdot)$ by (2.6).

(b) For given $g_j(\cdot)$, search β using an algorithm described in section 2.2.2.

Step 3: For $k = 1, \dots, q$, calculate $\text{GCV}(h_k)$ with β equal to its estimated value, where $\text{GCV}(\cdot)$ is defined in section 2.2.3. Let \hat{a}_1 and \hat{a}_2 be the minimiser of $\sum_{k=1}^q \{\text{GCV}(h_k) - a_0 - a_1 h_k^4 - a_2/(n h_k)\}^2$. Define $\hat{h} = \{\hat{a}_2/(4n\hat{a}_1)\}^{1/5}$ if \hat{a}_1 and \hat{a}_2 are positive, and $\hat{h} = \text{argmin}_{h_k} \text{GCV}(h_k)$ otherwise.

Step 4: For $h = \hat{h}$ selected in Step 3, repeat (a) and (b) in Step 2 until two successive values of $R(\beta)$ differ insignificantly.

Step 5: For $\beta = \hat{\beta}$ selected in Step 4, apply the stepwise deletion of section 2.2.4 to select significant variables X_{tj} for each fixed z .

Some additional remarks are now in order.

Remark 2. (i) The standardisation in Step 1 also ensures that the sample mean of $\{\beta^T \mathbf{X}_t\}$ is 0 and the sample variance is 1 for any unit vector β . This effectively re-writes model (2.3) as

$$\sum_{j=0}^d g_j \left(\beta^T \hat{\Sigma}^{-1/2} (\mathbf{x} - \hat{\mu}) \right) x_j,$$

where $\hat{\mu}$ and $\hat{\Sigma}$ are the sample mean and sample variance, respectively. In the numerical examples in section 4, we report $\hat{\Sigma}^{-1/2} \hat{\beta} / \|\hat{\Sigma}^{-1/2} \hat{\beta}\|$ as the estimated value of β defined in (2.3).

(ii) We may let $w(z) = I(|z| \leq 2 + \delta)$ for some small $\delta \geq 0$. To further speed up the computation, we estimate functions $g_j(\cdot)$ in Step 3 on 101 regular grids in the interval $[-1.5, 1.5]$ first, and then estimate the values of the functions on this interval by linear interpolation. Finally in Step 4, we estimate functions $g_j(\cdot)$ on the interval $[-2, 2]$.

(iii) With the Epanechnikov kernel, we let $q = 15$ and $h_k = 0.2 \times 1.2^{k-1}$ in Step 3. The specified values for h practically cover the range of 0.2 to 2.57 times of the standard deviation of the data. If we use the Gaussian kernel, we may select the range of the bandwidth between 0.1 and 1.5 times of the standard deviation.

(iv) To further stabilise the search for β in Step 2(b), we replace an estimate of $g_j(\cdot)$ on a grid point by a weighted average on its 5 nearest neighbours with weights $\{1/2, 1/6, 1/6, 1/12, 1/12\}$. The edge points are adjusted accordingly.

(v) In searching for β in terms of the one-step iterative algorithm, we estimate the derivatives of $g_j(\cdot)$ based on their adjusted estimates on the grid points as follows:

$$\hat{g}_j(z) = \{\hat{g}_j(z_1) - \hat{g}_j(z_2)\} / (z_1 - z_2), \quad j = 0, \dots, d,$$

$$\hat{\hat{g}}_j(z) = \{\hat{g}_j(z_1) - 2\hat{g}_j(z_2) + \hat{g}_j(z_3)\} / (z_1 - z_2)^2, \quad j = 0, \dots, d,$$

where $z_1 > z_2 > z_3$ are three nearest neighbours of z among the 101 regular grid points (see (ii) above). Equation (2.8) should be iterated a few times instead of just once.

(vi) Although the proposed algorithm works well with the examples reported in section 4, its convergence requires further research. In practice, we may detect if an estimated β is likely to be the global minimum by using multiple initial values.

3 Varying-coefficient linear models with two indices

In this section, we consider varying-coefficient models with two indices but one of them known. We assume knowing one index in order to keep computation practically feasible.

Let Y and V be two random variables, and \mathbf{X} be a $d \times 1$ random vector. We use V to denote the known index, which could be a (known) linear combination of \mathbf{X} . The goal is to approximate

the conditional expectation $G(\mathbf{x}, v) = E(Y|\mathbf{X} = \mathbf{x}, V = v)$, in the mean square sense (see (2.1)), by a function of the form

$$g(\mathbf{x}, v) = \sum_{j=0}^{d-1} g_j(\beta^T \mathbf{x}, v) x_j, \quad (3.1)$$

where $\beta = (\beta_1, \dots, \beta_d)^T$ is a $d \times 1$ unknown unit vector. Similar to Theorem 1(ii), it may be proved that under some mild conditions on $g(\mathbf{x}, v)$, the expression on the RHS of (3.1) is unique if the first non-zero β_k is positive and $\beta_d \neq 0$. Let $\{(\mathbf{X}_t, V_t, Y_t); 1 \leq t \leq n\}$ be observations from a strictly stationary process, and (\mathbf{X}_t, V_t, Y_t) has the same distribution as (\mathbf{X}, V, Y) .

The estimation for $g(\mathbf{x}, v)$ can be carried out in a similar manner as in one index case (see section 2.3). We outline the algorithm below briefly.

Step 1: Standardise the data $\{\mathbf{X}_t\}$ such that it has sample mean 0 and sample covariance matrix \mathbf{I}_d . Standardise the data $\{V_t\}$ such that V_t has sample mean 0 and sample variance 1. Specify an initial value of β .

Step 2: For each prescribed bandwidth value h_k , $k = 1, \dots, q$, repeat (a) and (b) below until two successive values of $R(\beta)$ defined in (3.2) differ insignificantly.

(a) For a given β , estimate $g_j(\cdot, \cdot)$ in terms of local linear regression.

(b) For given $g_j(\cdot, \cdot)$, search for β using a one-step iteration algorithm.

Step 3: For $k = 1, \dots, q$, calculate $\text{GCV}(h_k)$ with β equal to its estimated value, where $\text{GCV}(\cdot)$ is defined as in section 2.2.3. Let \hat{a}_1 and \hat{a}_2 be the minimiser of $\sum_{k=1}^q \{\text{GCV}(h_k) - a_0 - a_1 h_k^4 - a_2/(n h_k^2)\}^2$. Define $\hat{h} \equiv \{\hat{a}_2/(2 n \hat{a}_1)\}^{1/6}$.

Step 4: For $h = \hat{h}$ selected in Step 3, repeat (a) and (b) in Step 2 until two successive values of $R(\beta)$ differ insignificantly.

Step 5: For $\beta = \hat{\beta}$ from Step 4, select local significant variables for each fixed (z, v) .

Remark 3. (i) In Step 2(a) above, The local linear regression estimation leads to the problem of minimising the sum

$$\sum_{t=1}^n \left[Y_t - \sum_{j=0}^{d-1} \{a_j + b_j(\beta^T \mathbf{X}_t - z) + c_j(V_t - v)\} X_{tj} \right]^2 K_h(\beta^T \mathbf{X}_t - z, V_t - v) w(\beta^T \mathbf{X}_t, V_t),$$

where $K_h(z, v) = h^{-2} K(z/h, v/h)$, $K(\cdot, \cdot)$ is a kernel function on \mathbb{R}^2 , and $w(\cdot, \cdot)$ is a bounded weight function with a bounded support in \mathbb{R}^2 . We use a common bandwidth h for simplicity. The derived estimators are $\hat{g}_j(z, v) = \hat{a}_j$, $\hat{g}_{j,z}(z, v) = \hat{b}_j$ and $\hat{g}_{j,v}(z, v) = \hat{c}_j$ for $j = 0, \dots, d-1$, where $\dot{g}_{j,z}(z, v) = \partial g_j(z, v)/\partial z$ and $\dot{g}_{j,v}(z, v) = \partial g_j(z, v)/\partial v$.

(ii) In Step 2(b), we search for β which minimises the function

$$R(\beta) = \frac{1}{n} \sum_{t=1}^n \left\{ Y_t - \sum_{j=0}^{d-1} g_j(\beta^T \mathbf{X}_t, V_t) X_{tj} \right\}^2 w(\beta^T \mathbf{X}_t, V_t). \quad (3.2)$$

A one-step iterative algorithm may be constructed for this purpose in a similar manner as in section 2.2.2. The required estimates for the second derivatives of $g_j(z, v)$ may be obtained via a partially local quadratic regression.

(iii) In Step 3, the estimated $g(\mathbf{x}, v)$ is linear in $\{Y_t\}$ (for a given β). Thus, the generalized cross-validation method outlined in section 2.2.3 is still applicable.

(iv) Locally around given indices $\beta^T \mathbf{x}$ and v , (3.1) is approximately a linear model. Thus, the local variable selection technique outlined in section 2.2.4 is still applicable in Step 5.

4 Numerical properties

We use the Epanechnikov kernel in our calculation. The one-step iterative algorithm described in section 2.2.2 is used to estimate the index β , in which we iterate the ridge version of equation (2.8) two to four times. We stop the search in Step 2 when either the two successive values of $R(\beta)$ differ less than 0.001, or the number of replications of (a) and (b) in Step 2 exceeds 30. Setting initially the ridge parameter $q_n = 0.001 n^{-1/2}$, we keep doubling its value until $\ddot{R}_r(\cdot)$ is no longer ill-conditioned with respect to the precision of computers.

4.1 Simulation

We demonstrate the finite-sample performance of the varying-coefficient model with one index in Example 1, and with two indices in Example 2. We use the absolute inner product $|\beta^T \hat{\beta}|$ to measure the goodness of the estimated direction $\hat{\beta}$. Their inner product represents the cosine of the angles between the two directions. For Example 1, we evaluate the performance of the estimator in terms of the mean absolute deviation error

$$\mathcal{E}_{\text{MAD}} = \frac{1}{101 d} \sum_{j=0}^{d-1} \sum_{k=1}^{101} |\hat{g}_j(z_k) - g_j(z_k)|,$$

where z_k , $k = 1, \dots, 101$, are the regular grid points on $[-2, 2]$ after the standardisation. For Example 2, \mathcal{E}_{MAD} is calculated on the observed values instead.

Example 1. Consider the regression model

$$\begin{aligned} Y_t &= 3 \exp\{-Z_t^2\} + 0.8 Z_t X_{t1} + 1.5 \sin(\pi Z_t) X_{t3} + \varepsilon_t, \\ \text{with } Z_t &= \frac{1}{3}(X_{t1} + 2X_{t2} + 2X_{t4}), \end{aligned}$$

where $\mathbf{X}_t \equiv (X_{t1}, \dots, X_{t4})^T$, for $t \geq 1$, are independent random vectors uniformly distributed on $[-1, 1]^4$, and ε_t are independent $N(0, 1)$ random variables. The regression function in the above model is of form (2.3) with $d = 4$, $\beta = \frac{1}{3}(1, 2, 0, 2)^T$, and

$$g_0(z) = 3e^{-z^2}, \quad g_1(z) = 0.8z, \quad g_2(z) \equiv 0, \quad \text{and } g_3(z) = 1.5 \sin(\pi z).$$

We conduct two simulations with sample size 200 and 400 respectively, each with 200 replications. CPU time for each replication with sample size 400 is about 18 seconds on a Pentium II 350MHz

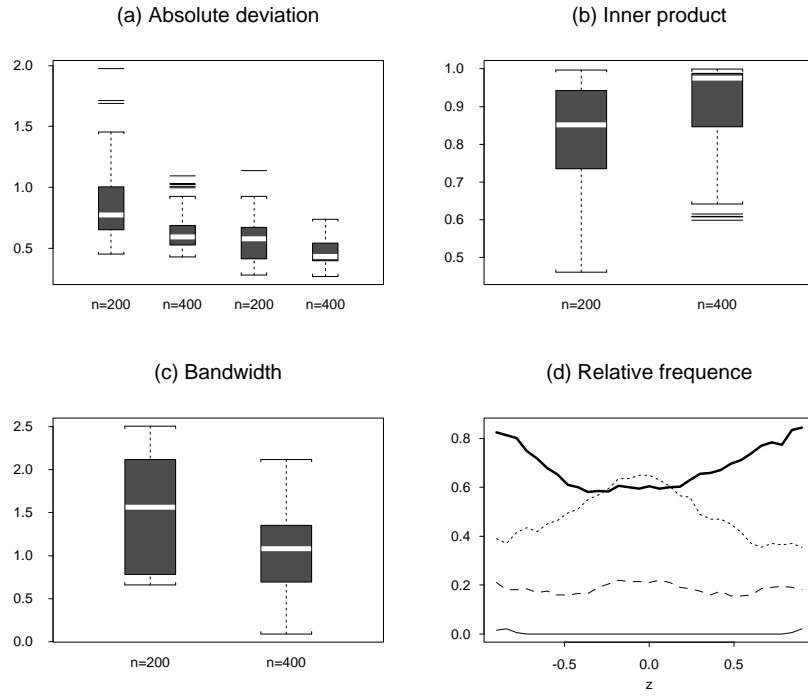


Figure 1: *Simulation results for Example 1. (a) Boxplots of \mathcal{E}_{MAD} . The two panels on the left are based on $\hat{\beta}$, and the two panels on the right are based on the true β . (b) Boxplots of $|\beta^T \hat{\beta}|$. (c) Boxplots of selected bandwidths. (d) Plots of the relative frequencies for deletion of locally insignificant terms at z against z : thin solid line — for the intercept; dotted line — for X_{t1} , thick solid line — for X_{t2} , and dashed line — for X_{t3} .*

PC (Linux). The results are summarised in Fig.1. Fig.1(a) displays boxplots of the mean absolute deviation errors. We also plot those errors obtained using the true direction β . The deficiency due to unknown β decreases when the sample size increases. Fig.1(b) shows that the estimator $\hat{\beta}$ derived from the one-step iterative algorithm is close to the true β with high frequency. The average iteration time in search for β is 14.43 for $n = 400$ and 18.25 for $n = 200$. Most outliers in Fig.1(a) and Fig.1(b) correspond to the cases where the search for β does not converge within 30 iterations. Fig.1(c) indicates that the proposed bandwidth selector is stable. We also apply the method in section 2.2.4 to choose the local significant variables at the 31 regular grid points in the range from -1.5 to 1.5 times the standard deviations of $\beta^T \mathbf{X}$. The relative frequencies of deletion are depicted in Fig.1(d). There is overwhelming evidence to include the ‘intercept’ $g_0(z) = 3e^{-z^2}$ in the model for all the values of z . In contrast, we tend to delete most often the term X_{t2} which has ‘coefficient’ $g_2(z) \equiv 0$. There is strong evidence to keep the term X_{t3} in the model. Note that the term X_{t1} is less significant, as the magnitude of its ‘coefficient’ $g_1(z) = 0.8z$ is smaller than that of both $g_0(z)$ and $g_3(z)$.

Fig.2 presents three typical examples of the estimated coefficient functions with the sample size $n = 400$. The curves are plotted on the range from -1.5 to 1.5 times the standard deviation of $\beta^T \mathbf{X}$. The three examples are selected with the corresponding \mathcal{E}_{MAD} at the first quartile, the

median and the third quartile respectively among the 200 replications. For the example with \mathcal{E}_{MAD} at the median, we also plot the estimated functions obtained using the true index β . For that example, $\hat{\beta}^T \beta = 0.946$. The deficiency due to unknown β is almost negligible. Note that the biases of the estimators for the coefficient functions $g_0(\cdot)$, $g_1(\cdot)$ and $g_2(\cdot)$ are large near to boundaries. We believe that this is due to the collinearity of functions g_i and small effective local sample size near the tails. Nevertheless, there seems no evidence that the problem have distorted the estimation for the target function $g(\mathbf{x})$.

We have also repeated the above exercise with $\varepsilon_t \sim N(0, \sigma^2)$ for different values of σ^2 . While the results are in the same pattern as the above, the estimation for both coefficient functions and the direction β is more accurate for the models with smaller noise.

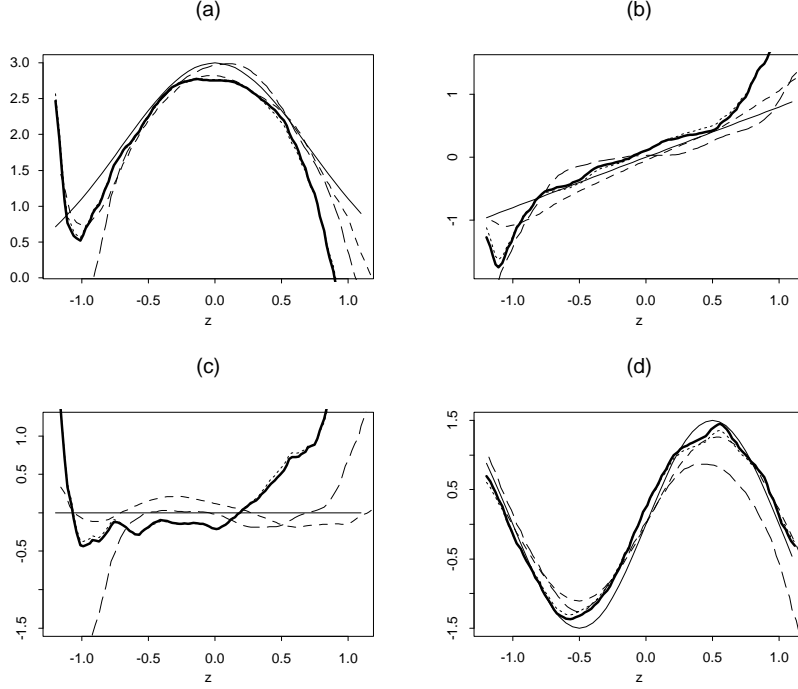


Figure 2: *Simulation results for Example 1 ($n = 400$). Estimated coefficient functions with \mathcal{E}_{MAD} at first quartile (short dashed curve), at the third quartile (long dashed curve), and at median (thick curve) with the true β (dotted curve), together with true functions (thin solid curve). (a) $g_0(z) = 3e^{-z^2}$; (b) $g_1(z) = 0.8z$; (c) $g_2(z) = 0$; (d) $g_3(z) = 1.5 \sin(\pi z)$.*

Example 2. We consider the regression model

$$Y_t = 3 \exp(-Z_t^2 + X_{t1}) + (Z_t + X_{t1}^2)X_{t1} - \log(Z_t^2 + X_{t1}^2)X_{t2} + 1.5 \sin(\pi Z_t + X_{t1})X_{t3} + \varepsilon_t,$$

$$\text{with } Z_t = \frac{1}{2}(X_{t1} + X_{t2} + X_{t3} + X_{t4}),$$

where $\{X_{t1}, \dots, X_{t4}\}$ and $\{\varepsilon_t\}$ are the same as in Example 1. Obviously, the regression function in the above model is of the form (3.1) with $d = 4$, $\beta = \frac{1}{2}(1, 1, 1, 1)^T$, $V_t = X_{t1}$ and

$$g_0(z, v) = 3e^{-z^2+v}, \quad g_1(z, v) = z + v^2, \quad g_2(z, v) = -\log(z^2 + v^2), \quad g_3(z, v) = 1.5 \sin(\pi z + v).$$

We conduct three simulations with sample size 200, 400 and 600 respectively, each with 100 replications. CPU time for each realisation, in a Sun Ultra-10 300MHz Workstation, is about 18 seconds for $n = 200$, 80 seconds for $n = 400$ and 190 seconds for $n = 600$. Fig.3(a) shows that the mean absolute deviation error \mathcal{E}_{MAD} decreases when n increases. For the sake of comparison, we also present \mathcal{E}_{MAD} based on the true β . Fig.3(b) displays the boxplots of the absolute inner product $|\beta^T \hat{\beta}|$, which indicates that the one-step iteration algorithm works reasonably well. The boxplots of bandwidths selected by the GCV-method are depicted in Fig.3(c).

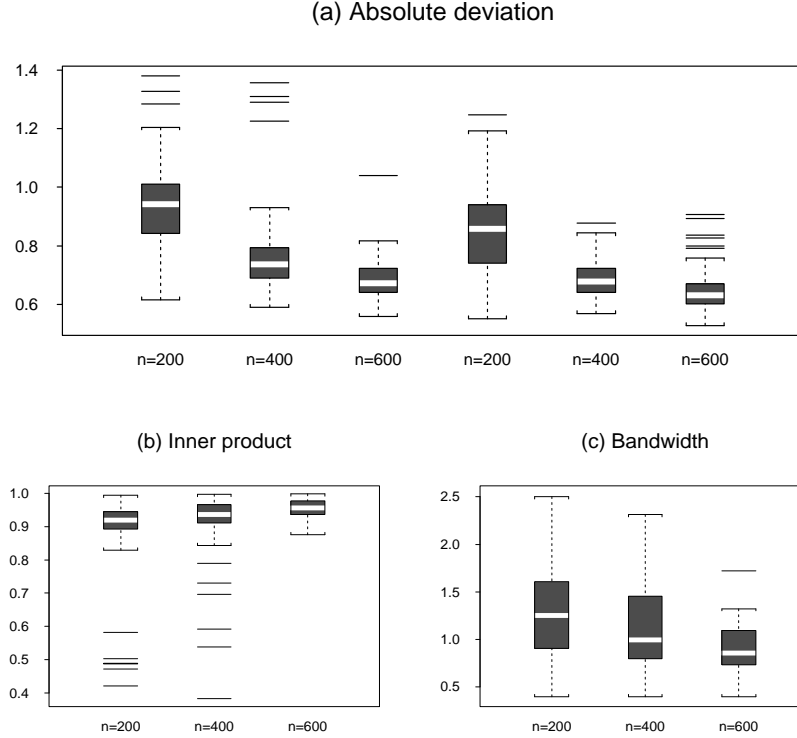


Figure 3: *Simulation results for Example 2. Boxplots of (a) \mathcal{E}_{MAD} , (b) $|\beta^T \hat{\beta}|$, and (c) the selected bandwidths. The three panels on the left in (a) are based on $\hat{\beta}$, and the three panels on the right are based on the true β .*

4.2 Real data examples

Example 3. The annual numbers of muskrats and mink caught over 82 trapping regions have been recently extracted from the records compiled by the Hudson Bay Company on fur sales at auction in 1925-1949. Fig.4(a) indicates the 82 posts where furs were collected. Fig.4(b) plots the time series of the mink and the muskrat (on the natural logarithmic scale) from 8 randomly selected posts. There exists a clear synchrony between the fluctuations of the two species with a delay of one or two years; indicating the food chain interaction between prey (*i.e.* muskrat) and predator (*i.e.* mink); see Errington (1963). A simple biological model for food chain interaction proposed

(a) The 82 trapping posts for the mink and the muskrat in Canada

(Put Fig.4(a) here)

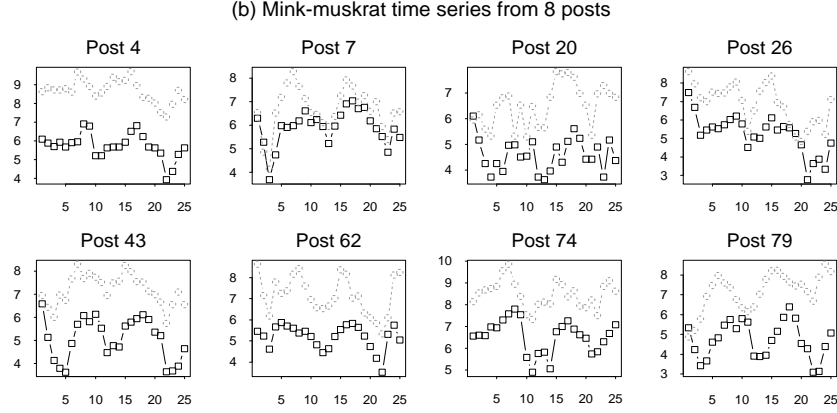


Figure 4: (a) Map of 82 posts for the mink and the muskrat in Canada in 1925 – 1949. (b) Time series plots of the mink and the muskrat data from 8 randomly selected posts. Solid lines — mink; dashed lines — muskrats.

by May (1981) and Stenseth *et al.* (1997) is of the form

$$\begin{cases} X_{t+1} - X_t &= a_0(\theta_t) - a_1(\theta_t)X_t - a_2(\theta_t)Y_t, \\ Y_{t+1} - Y_t &= b_0(\theta_t) - b_1(\theta_t)Y_t + b_2(\theta_t)X_t, \end{cases} \quad (4.1)$$

where X_t and Y_t denote the population abundances, on a natural logarithmic scale, of prey (muskrat) and predator (mink) respectively at time t , $a_i(\cdot)$ and $b_i(\cdot)$ are non-negative functions, and θ_t is an indicator representing *the regime effect* at time t , which is determined by X_t and/or Y_t . The term ‘regime effect’ collectively refers to the nonlinear effect due to, among others, the different hunting/escaping behaviour and reproduction rates of animals at different stages of population fluctuation (Stenseth *et al.*, 1999). In fact, $a_1(\theta_t)$ and $b_1(\theta_t)$ reflect within species regulation whereas $a_2(\theta_t)$ and $b_2(\theta_t)$ reflect the food chain interaction between the two species.

Model (4.1), with added random noise, would be in the form of varying-coefficient linear models if we let θ_t be a linear combination X_t and Y_t . However each mink and muskrat time series has only 25 points, which is too short for fitting such a nonlinear model. Based on some statistical tests on the common structure for each pair among those 82 posts, Yao *et al.* (2000) suggested a grouping with three clusters: the eastern area consisting of post 10, post 67 and the other six posts on its right in Fig.4(a); the western area consisting of the 30 posts on the left in Fig.4(a) (*i.e.* post 17 and those on its left); and the central area consisting of the remaining 43 posts in the middle. Since some data are missing at post 15, we exclude it from our analysis. The sample size for eastern, central and western areas are therefore 207, 989 and 667 respectively. With the new technique proposed in this paper, we fit the pooled data for each of the three areas with the model

$$\begin{cases} X_{t+1} &= f_0(Z_t) + f_1(Z_t)Y_{t-1} + f_2(Z_t)Y_t + f_3(Z_t)X_{t-1} + \varepsilon_{1,t+1}, \\ Y_{t+1} &= g_0(Z_t) + g_1(Z_t)Y_{t-1} + g_2(Z_t)Y_t + g_3(Z_t)X_{t-1} + \varepsilon_{2,t+1}, \end{cases} \quad (4.2)$$

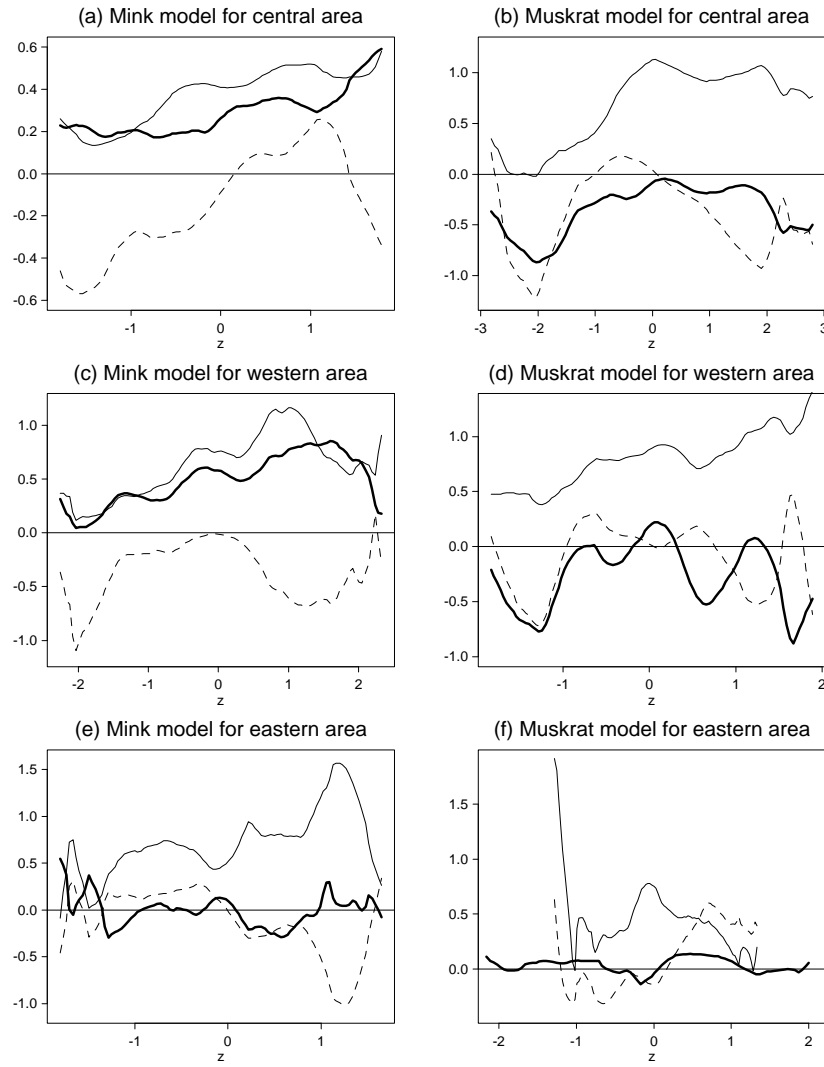


Figure 5: *Estimated coefficient functions for Canadian mink-muskrat data. (a), (c) & (d): thick solid lines — $g_x(\cdot)$; solid lines — $g_y(\cdot)$; dashed lines — $g_0(\cdot)$. (b), (d) & (f): thick solid lines — $f_y(\cdot)$; solid lines — $f_x(\cdot)$; dashed lines — $f_0(\cdot)$.*

where $Z_t = \beta_1 Y_{t-1} + \beta_2 Y_t + \beta_3 X_{t-1} + \beta_4 X_t$ with $\beta \equiv (\beta_1, \beta_2, \beta_3, \beta_4)^T$ selected by data. Comparing with (4.1), we include further lagged values X_{t-1} and Y_{t-1} into the above model. To eliminate the effect of different sampling weights in different regions and for different species, we first standardised mink and muskrat series separately for each post. We apply the local deletion technique presented in section 2.2.4 to detect *local* redundant variables at 31 regular grid points over the range from -1.5 to 1.5 times the standard deviation of Z_t . As a primary attempt, we also select the *global* model based on the local deletion. A more direct approach would be, for example, based on the generalized likelihood ratio tests of Fan, Zhang and Zhang (2001). We denote by R_{MSE} the ratio of the mean squared errors from the fitted model over the sample variance of the variable to be fitted.

First, we use the second equation of (4.2) to model mink population dynamics in the central area.

The selected β is $(0.424, 0.320, 0.432, 0.733)^T$, the selected bandwidth is 0.415, and $R_{MSE} = 0.449$. The local variable selection indicates that X_{t-1} is the least significant over all, for it is significant at only 7 out of 31 grid points. By leaving it out, we reduce to the model

$$Y_{t+1} = g_0(Z_t) + g_y(Z_t)Y_t + g_x(Z_t)X_t + \varepsilon_{2,t+1}, \quad (4.3)$$

where $Z_t = \beta_1 Y_t + \beta_2 X_t + \beta_3 Y_{t-1}$. Our algorithm selects

$$Z_t = (0.540Y_t - 0.634Y_{t-1}) + 0.553X_t, \quad (4.4)$$

which suggests that the nonlinearity is dominated by the growth rate of mink (i.e. $Y_t - Y_{t-1}$) and the population of muskrat (i.e. X_t) in the previous year. The estimated coefficient functions are plotted in Fig.5(a). The coefficient function $g_x(\cdot)$ is positive, which reflects the fact that a large muskrat population will facilitate the growth of the mink population. The coefficient function $g_y(\cdot)$ is also positive, which reflects a natural reproduction process of mink population. Both $g_y(\cdot)$ and $g_x(\cdot)$ are approximately increasing with respect to the sum of growth rate of mink and population of muskrat; see (4.4). All the terms in model (4.3) are significant in most places; the number of significant grid points for ‘intercept’, Y_t and X_t are 21, 31 and 26 (out of 31 in total). The selected bandwidth is 0.597 and $R_{MSE} = 0.461$.

Starting with the first equation of (4.2), the fitted model for the muskrat dynamics in the central area is

$$X_{t+1} = f_0(Z_t) + f_y(Z_t)Y_t + f_x(Z_t)X_t + \varepsilon_{1,t+1} \quad (4.5)$$

with $Z_t = 0.542Y_t + 0.720X_t + 0.435X_{t-1}$, $\hat{h} = 0.498$ and $R_{MSE} = 0.559$. The estimated coefficient functions are plotted in Fig.5(b). The coefficient function $f_y(\cdot)$ is always negative, which reflects the fact that mink is the key predator of muskrat in this core of the boreal forest in Canada. The coefficient $f_x(\cdot)$ is positive, as expected.

We repeat the above exercise for pooled data in the western area, resulting in similar results. In fact, model (4.3) appears appropriate for mink dynamics with $Z_t = 0.469Y_t + 0.723X_t + 0.507Y_{t-1}$, $R_{MSE} = 0.446$, $\hat{h} = 0.415$, and the estimated coefficient functions plotted in Fig.5(c). Model (4.5) appears appropriate for muskrat dynamics with $Z_t = 0.419Y_t + 0.708X_t + 0.569X_{t-1}$, $\hat{h} = 0.415$, $R_{MSE} = 0.416$, and the estimated coefficient functions plotted in Fig.5(d).

Fitting the data in the eastern area leads to drastically different results from the above. The fitting for the mink dynamics with model (4.3) gives $Z_t = 0.173Y_t - 0.394X_t + 0.901Y_{t-1}$, $\hat{h} = 0.597$ and $R_{MSE} = 0.681$. Out of 31 grid points, the ‘intercept’, Y_t and X_t are significant at 15, 31 and 4 points respectively. There is clear auto-dependence in the mink series $\{Y_t\}$ while the muskrat data $\{X_t\}$ carry little information about minks. The estimated coefficients, depicted in Fig.5(e), reinforce the above observation. The fitting of the muskrat dynamics shows again that there seems little interaction between mink and muskrat in this area. For example, the term Y_t in model (4.5) is not significant at all the 31 grid points. The estimated coefficient function $f_y(\cdot)$ is plotted as the

thick curve in Fig.5(f), which is always close to 0. We fit the data with a further simplified model

$$X_{t+1} = f_0(Z_t) + f_x(Z_t)X_t + \varepsilon_{1,t+1},$$

resulting $Z_t = 0.667X_t - 0.745X_{t-1}$, $\hat{h} = 0.498$ and $R_{MSE} = 0.584$. The estimated coefficient functions are superimposed on Fig.5(f). Note the different ranges of z -values are due to different Z_t 's used in the above model and model (4.5).

In summary, we have facilitated the data analysis of the biological food chain interaction model of Stenseth *et al.* (1997) by portraying the nonlinearity through varying-coefficient linear forms. The selection of the index in our algorithm is equivalent in this context to the selection of the regime effect indicator, which in itself is of biological interest. The numerical results indicate that there is a strong evidence of the predator-prey interactions between the minks and the muskrats in the central and western areas. However, no evidence for such an interaction exists in the eastern area. In light of what is known in the eastern area, this is not surprising. There is a larger array of prey-species for the mink to feed on, making it less dependent on muskrat (see Elton, 1942).

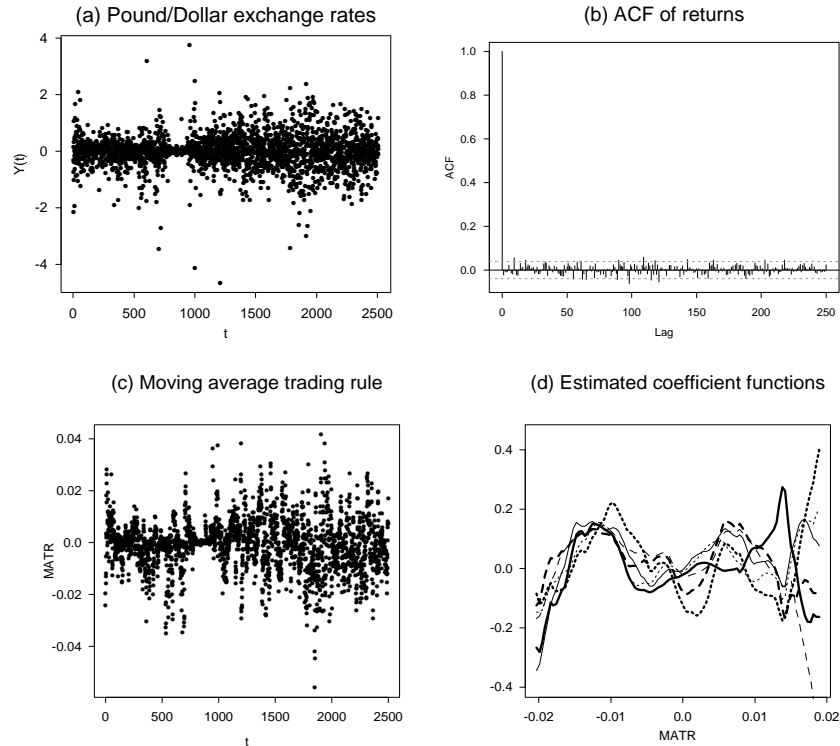


Figure 6: (a) Pound/Dollar exchange rate return series $\{Y_t\}$. (b) Autoregressive function of $\{Y_t\}$. (c) Plot of $\{U_t = Y_t/(\sum_{i=0}^9 Y_{t-i}/10)\}$. (d) Estimated coefficient functions of model (4.6) with $Z_t = U_{t-1}$ and $m = 5$. Thick solid lines — $g_0(\cdot)$, thick dotted lines — $g_1(\cdot)$, thick dashed lines — $g_2(\cdot)$, solid lines — $g_3(\cdot)$, dotted lines — $g_4(\cdot)$, dashed lines — $g_5(\cdot)$.

Example 4. This example concerns the daily closing bid prices of the pound sterling in terms of US dollar from 2 January 1974 to 30 December 1983, which forms a time series of length 2510. The previous analysis of this ‘particularly difficult’ data set can be found in Gallant, Hsieh and

Tauchen (1991) and the references within. Let X_t be the exchange rate on the t -th day. We model the return series $Y_t = 100 \log(X_t/X_{t-1})$, plotted in Fig.6(a), using the techniques developed in this paper. Typically the classical financial theory would treat $\{Y_t\}$ as a martingale difference process. Therefore Y_t would be unpredictable. Fig.6(b) shows that there exists almost no significant autocorrelation in $\{Y_t\}$.

First, we approximate the conditional expectation of Y_t (given its past) by

$$g_0(Z_t) + \sum_{i=1}^m g_i(Z_t)Y_{t-i}, \quad (4.6)$$

where $Z_t = \beta_1 Y_{t-1} + \beta_2 Y_{t-2} + \beta_3 X_{t-1} + \beta_4 U_{t-1}$, and $U_{t-1} = X_{t-1} \left\{ L^{-1} \sum_{j=1}^L X_{t-j} \right\}^{-1} - 1$. The variable U_{t-1} defines the *moving average technical trading rule* (MATR) in finance, and $U_{t-1} + 1$ is the ratio of exchange rate at the time $t - 1$ to the average rate over past period of length L . The MATR signals 1 (the position to *buy* sterling) when $U_{t-1} > 0$, and -1 (the position to *sell* sterling) when $U_{t-1} < 0$. For detailed discussion of the MATR, we refer to the papers by LeBaron (1997, 1999) and Hong and Lee (1999). We use the first 2410 sample points for estimation and last 100 points for post-sample forecasting. We evaluate the post-sample forecast by the *mean trading return* defined as

$$\text{MTR} = \frac{1}{100} \sum_{t=1}^{100} S_{2410+t-1} Y_{2410+t},$$

where S_t is a signal function taking values -1 , 0 and 1 . The mean trading return measures the real profits in a financial market, ignoring interest differentials and transaction costs (for the sake of simplicity). It is more relevant than the conventional mean squared predictive errors or average absolute predictive errors for evaluating the performance of forecasting for market movements; see Hong and Lee (1999). Under this criterion, we need to predict the direction of market movement rather than its quantity. For the MATR, the mean trading return is defined as

$$\text{MTR}_{\text{MA}} = \frac{1}{100} \sum_{t=1}^{100} \{I(U_{2410+t-1} > 0) - I(U_{2410+t-1} < 0)\} Y_{2410+t}.$$

Let \hat{Y}_t be defined as the estimated function given in (4.6). The mean trading return for the forecasting based on our varying-coefficient modeling is defined as

$$\text{MTR}_{\text{VC}} = \frac{1}{100} \sum_{t=1}^{100} \{I(\hat{Y}_{2410+t} > 0) - I(\hat{Y}_{2410+t} < 0)\} Y_{2410+t}.$$

On the other hand, ideally we would buy at time $t - 1$ when $Y_t > 0$ and sell when $Y_t < 0$. The mean trading return for this ‘ideal’ strategy is

$$\text{MTR}_{\text{ideal}} = \frac{1}{100} \sum_{t=1}^{100} |Y_{2410+t}|,$$

which serves as a benchmark for assessing forecasting procedures. For example, for this particular data set, $\text{MTR}_{\text{MA}}/\text{MTR}_{\text{ideal}} = 12.58\%$ if we let $L = 10$.

Now we are ready to proceed. First, we let $m = 5$ and $L = 10$ in (4.6), *i.e.* we use one week data in the past as the ‘regressors’ in the model and define the MATR by comparing with the average rate in last two weeks. The selected β is $(0.0068, 0.0077, 0.0198, 0.9998)^T$ which suggests that U_t plays an important role in the underlying nonlinear dynamics. The ratio of the MSE of the fitted model to the sample variance of $\{Y_t\}$ is 93.67%, which reflects the presence of high level ‘noise’ in the financial data. The selected bandwidth is 0.24. The ratio $\text{MTR}_{\text{VC}}/\text{MTR}_{\text{ideal}} = 5.53\%$. The predictability is much lower than that of the MATR. If we include rates in last two weeks as regressors in the model (*i.e.* $m = 10$ in (4.6)), the ratio $\text{MTR}_{\text{VC}}/\text{MTR}_{\text{ideal}}$ increases to 7.26% which is still distance away from $\text{MTR}_{\text{MA}}/\text{MTR}_{\text{ideal}}$, while the ratio of the MSE of the fitted model to the sample variance of $\{Y_t\}$ is 87.96%. The selected bandwidth is still 0.24, and $\hat{\beta} = (0.0020, 0.0052, 0.0129, 0.9999)^T$.

The above calculations (also others not reported here) seem to suggest that U_t could be a dominant component in the selected index. This leads us to use model (4.6) with fixed $Z_t = U_{t-1}$, which is actually the approach adopted by Hong and Lee (1999). For $m = 5$, the fitting to the data used in estimation becomes worse; the ratio of the MSE of the fitted model to the sample variance of $\{Y_t\}$ is 97.39%. But it provides a better post-sample forecasting; $\text{MTR}_{\text{VC}}/\text{MTR}_{\text{ideal}}$ is 23.76%. The selected bandwidth is 0.24. The plots of estimated coefficient functions indicate a possible under-smoothing. By increasing the bandwidth to 0.40, $\text{MTR}_{\text{VC}}/\text{MTR}_{\text{ideal}}$ is 31.35%. The estimated coefficient functions are plotted in Fig.6(d). The rate of correct predictions for the direction of market movement (*i.e.* sign of Y_t) is 50% for the MATR, and 53% and 58% for the varying-coefficient model with bandwidth 0.24 and 0.40 respectively.

A word of caution: We should not take for granted the above improvement in forecasting from using U_t as the index. Hong and Lee (1999) conducted empirical studies of this approach with several financial data sets with only partial success. In fact, for this particular data set, model (4.6) with $Z_t = U_t$ and $m = 10$ gives a negative value of MTR_{VC} . Note that the ‘super-dominating’ position of U_t in the selected smoothing variable $\hat{\beta}^T \mathbf{X}_t$ is partially due to the scaling difference between U_t and (Y_t, X_t) ; see also Fig.6(a) and Fig.6(c). In fact, if we standardise U_t , Y_t and X_t separately beforehand, the resulting $\hat{\beta}$ is $(0.59, -0.52, 0.07, 0.62)^T$ when $m = 5$, which is dominated by U_{t-1} and the contrast between Y_{t-1} and Y_{t-2} . ($\text{MTR}_{\text{VC}}/\text{MTR}_{\text{ideal}} = 1.42\%$. The ratio of MSE of the fitted model to the sample variance of Y_t is 96.90%.) By doing this, we effectively use a different class of models to *approximate* the unknown conditional expectation of Y_t ; see Remark 2(i). Finally, we remark that a different modeling approach should be adopted if our primary target is to maximise the mean trading return, which is obviously beyond the scope of this paper.

Appendix: Proof of Theorem 1

We use the same notation as in section 2.

Proof of Theorem 1(i). It follows from the ordinary least-squares theory that there exists a

minimal value of

$$E \left[\{Y - f(X)\}^2 \mid \boldsymbol{\alpha}^T \mathbf{X} = z \right]$$

over the class of functions of the form $f(\mathbf{x}) = \sum_{i=0}^d f_i(\boldsymbol{\alpha}^T \mathbf{x}) x_i$ with all f_i measurable. Let $f_0^*(z), \dots, f_{d-1}^*(z)$ be the minimiser. Then

$$\{f_1^*(z), \dots, f_d^*(z)\}^T = \left\{ \text{var}(\mathbf{X} \mid \boldsymbol{\alpha}^T \mathbf{X} = z) \right\}^- \text{cov}(\mathbf{X}, Y \mid \boldsymbol{\alpha}^T \mathbf{X} = z),$$

$$f_0^*(z) = E(Y \mid \boldsymbol{\alpha}^T \mathbf{X} = z) - \sum_{j=1}^d f_j^*(z) E(X_j \mid \boldsymbol{\alpha}^T \mathbf{X} = z).$$

In the above expression, A^- denotes a generalized inverse matrix of A for which $AA^-A = A$. It follows immediately from the least-squares theory that

$$E \left(\left\{ Y - f_0^*(z) - \sum_{j=1}^d f_j^*(z) X_j \right\}^2 \mid \boldsymbol{\alpha}^T \mathbf{X} = z \right) \leq \text{var}(Y \mid \boldsymbol{\alpha}^T \mathbf{X} = z)$$

Consequently,

$$R(\boldsymbol{\alpha}) \equiv E \left\{ Y - f_0^*(\boldsymbol{\alpha}^T \mathbf{X}) - \sum_{j=1}^d f_j^*(\boldsymbol{\alpha}^T \mathbf{X}) X_j \right\}^2$$

is bounded from the above by $\text{var}(Y)$, and continuous on the compact set $\{\boldsymbol{\alpha} \in R^d \mid \|\boldsymbol{\alpha}\| = 1\}$. Hence, there exists a $\boldsymbol{\beta}$ in the above set such that $R(\boldsymbol{\alpha})$ obtains its minimum at $\boldsymbol{\alpha} = \boldsymbol{\beta}$. Therefore, $g(\cdot)$ fulfilling (2.1) exists.

Theorem 1(ii) follows from the following two lemmas immediately. \square

Lemma A.1. Suppose that $F(\cdot) \not\equiv 0$ is a twice differentiable function defined on R^d , and

$$F(\mathbf{x}) = g_0(\boldsymbol{\beta}^T \mathbf{x}) + \sum_{j=1}^d g_j(\boldsymbol{\beta}^T \mathbf{x}) x_j \quad (\text{A.1})$$

$$= f_0(\boldsymbol{\alpha}^T \mathbf{x}) + \sum_{j=1}^d f_j(\boldsymbol{\alpha}^T \mathbf{x}) x_j, \quad (\text{A.2})$$

where $\boldsymbol{\alpha}, \boldsymbol{\beta}$ are non-zero and non-parallel vectors in R^d . Then $F(\mathbf{x}) = c_1 \boldsymbol{\alpha}^T \mathbf{x} \boldsymbol{\beta}^T \mathbf{x} + \boldsymbol{\gamma}^T \mathbf{x} + c_0$, where $\boldsymbol{\gamma} \in R^d$, $c_0, c_1 \in R$ are constants.

Proof. Without loss of the generality we assume $\boldsymbol{\beta} = (c, 0, \dots, 0)^T$. Then, it follows from (A.1) that $\partial^2 F(\mathbf{x}) / \partial x_i^2 = 0$ for $i = 2, \dots, d$. Write $\boldsymbol{\alpha}^T \mathbf{x} = z$. Choose $2 \leq i \leq d$ fixed for which $\alpha_i \neq 0$. Then from (A.2), we have that

$$\frac{\partial^2 F(\mathbf{x})}{\partial x_i^2} = \alpha_i^2 \ddot{f}_0(z) + \alpha_i^2 \sum_{j=1}^d \ddot{f}_j(z) x_j + 2\alpha_i \dot{f}_i(z) = 0,$$

namely,

$$\alpha_i \{ \alpha_i \ddot{f}_0(z) + z \ddot{f}_i(z) + 2\dot{f}_i(z) \} + \alpha_i^2 \sum_{j \neq i} \{ \ddot{f}_j(z) - \frac{\alpha_j}{\alpha_i} \ddot{f}_i(z) \} x_j = 0. \quad (\text{A.3})$$

Letting $x_j = 0$ for $j \neq i$ and $x_i = x / \alpha_i$ in the above equation, we have

$$\alpha_i \ddot{f}_0(x) + x \ddot{f}_i(x) + 2\dot{f}_i(x) = 0. \quad (\text{A.4})$$

Hence, (A.3) reduces to $\sum_{j \neq i} \{\ddot{f}_j(z) - \frac{\alpha_j}{\alpha_i} \ddot{f}_i(z)\} x_j = 0$, which leads to the equalities below if we let $x_k = x/\alpha_k$ and all other $x_j = 0$ for $k \neq i$ and $\alpha_k \neq 0$, or $x_k \neq 0$, $x_i = x/\alpha_i$ and all other $x_j = 0$ for $k \neq i$ and $\alpha_k = 0$:

$$\ddot{f}_k(x) = \ddot{f}_i(x) \frac{\alpha_k}{\alpha_i}, \quad 1 \leq k \leq d.$$

This implies that $f_k(z) = f_i(z) \alpha_k \alpha_i^{-1} + a_k z + b_k$ with $a_i = b_i = 0$. Substituting this into (A.2), we have

$$\begin{aligned} F(\mathbf{x}) &= f_0(\boldsymbol{\alpha}^T \mathbf{x}) + \alpha_i^{-1} f_i(\boldsymbol{\alpha}^T \mathbf{x}) \boldsymbol{\alpha}^T \mathbf{x} + \sum_{j \neq i} (a_j \boldsymbol{\alpha}^T \mathbf{x} + b_j) x_j \\ &\equiv f_0^*(\boldsymbol{\alpha}^T \mathbf{x}) + \sum_{j \neq i} (a_j \boldsymbol{\alpha}^T \mathbf{x} + b_j) x_j. \end{aligned}$$

Now, an application of the argument (A.4) to the last expression above shows that $f_0^*(z) = a_0 z + b_0$. Thus

$$F(\mathbf{x}) = a_0 \boldsymbol{\alpha}^T \mathbf{x} + b_0 + \sum_{j \neq i} (a_j \boldsymbol{\alpha}^T \mathbf{x} + b_j) x_j.$$

Now, $\partial^2 F(\mathbf{x}) / \partial x_i \partial x_j = a_j \alpha_i$ for any $j \geq 2$, which should be 0 according to (A.1) since $\boldsymbol{\beta} = (c, 0, \dots, 0)^T$. Hence, all a_j ($j \geq 2$) in the above expression are zero. This implies that

$$F(\mathbf{x}) = \boldsymbol{\gamma}^T \mathbf{x} + b_0 + a_1 x_1 \boldsymbol{\alpha}^T \mathbf{x} = \boldsymbol{\gamma}^T \mathbf{x} + b_0 + c^{-1} a_1 \boldsymbol{\beta}^T \mathbf{x} \boldsymbol{\alpha}^T \mathbf{x},$$

where $\boldsymbol{\gamma} = a_0 \boldsymbol{\alpha} + \mathbf{b}$, and $\mathbf{b} = (b_1, \dots, b_d)^T$. □

Lemma A.2. For any

$$F(\mathbf{x}) \equiv F(x_1, \dots, x_d) = f_0(\boldsymbol{\alpha}^T \mathbf{x}) + \sum_{j=1}^d f_j(\boldsymbol{\alpha}^T \mathbf{x}) x_j \neq 0,$$

where $\boldsymbol{\alpha} = (\alpha_1, \dots, \alpha_d)^T \in R^d$ and $\alpha_d \neq 0$, $F(\cdot)$ can be expressed as

$$F(\mathbf{x}) = g_0(\boldsymbol{\alpha}^T \mathbf{x}) + \sum_{j=1}^{d-1} g_j(\boldsymbol{\alpha}^T \mathbf{x}) x_j, \quad (\text{A.5})$$

where $g_0(\cdot), \dots, g_{d-1}(\cdot)$ are uniquely determined as follows:

$$g_0(z) = F(0, \dots, 0, z/\alpha_d), \quad (\text{A.6})$$

$$g_j(z) = F_j - g_0(z), \quad j = 1, \dots, d-1, \quad (\text{A.7})$$

where F_j denotes the value of F at $x_j = 1$, $x_d = (z - \alpha_j)/\alpha_d$ and $x_k = 0$ for all the other k .

Proof. Note that $x_d = \{\boldsymbol{\alpha}^T \mathbf{x} - \sum_{j=1}^{d-1} \alpha_j x_j\} / \alpha_d$. Define

$$g_0(z) = f_0(z) + \frac{1}{\alpha_d} f_d(z) z \quad \text{and} \quad g_j(z) = f_j(z) - \frac{\alpha_j}{\alpha_d} \quad \text{for } j = 1, \dots, d-1.$$

It is easy to see that (A.5) follows immediately. Let $x_1 = \dots = x_{d-1} = 0$ and $x_d = z/\alpha_d$ in (A.5), we obtain (A.6). Let $x_j = 1$, $x_d = (z - \alpha_j)/\alpha_d$ and $x_k = 0$ for all the other k , we obtain (A.7). The proof is completed. □

Acknowledgements

We thank Professors N.C. Stenseth and A.R. Gallant for making available Canadian mink-muskrat data and pound/dollar exchange data analysed in section 4.2. We are grateful to reviewers for very helpful comments.

References

- Bickel, P.J. (1975) One-step Huber estimates in linear models. *J. Amer. Statist. Assoc.*, **70**, 428–433.
- Cai, Z., Fan, J. and Li, R. (2000) Efficient estimation and inferences for varying-coefficient models. *J. Amer. Statist. Assoc.*, **95**, 888–902.
- Cai, Z., Fan, J. and Yao, Q. (2000) Functional-coefficient regression models for nonlinear time series models. *J. Amer. Statist. Assoc.*, **95**, 941–956.
- Carroll, R.J., Fan, J., Gijbels, I. and Wand, M.P. (1997) Generalized partially linear single-index models. *J. Amer. Statist. Assoc.*, **92**, 477–489.
- Carroll, R.J., Ruppert, D. and Welsh, A.H. (1998) Local estimating equations. *J. Amer. Statist. Assoc.*, **93**, 214–227.
- Chen, R. and Tsay, R.S. (1993) Functional-coefficient autoregressive models. *J. Amer. Statist. Assoc.*, **88**, 298–308.
- Cleveland, W.S., Grosse, E. and Shyu, W.M. (1992) Local regression models. In *Statistical Models in S* (ed. J.M. Chambers and T.J. Hastie), pp.309–376. Pacific Grove: Wadsworth & Brooks.
- Craven, P. and Wahba, G. (1979) Smoothing noisy data with spline functions: estimating the correct degree of smoothing by the method of generalized cross-validation. *Numerische Mathematik*, **31**, 377–403.
- Elton, C.S. (1942) *Voles, Mice and Lemmings*. Oxford: Clarendon Press.
- Errington, P.L. (1963) *Muskrat Populations*. Ames: Iowa State University Press.
- Fan, J. and Chen, J. (1999) One-step local quasi-likelihood estimation. *J. Roy. Statist. Soc. B*, **61**, 927–943.
- Fan, J. and Gijbels, I. (1996) *Local Polynomial Modeling and Its Applications*. London: Chapman and Hall.
- Fan, J., Zhang, C. and Zhang, J. (2001) Generalized likelihood ratio statistics and Wilks phenomenon. *Ann. Statist.*, **29**, 153–193.
- Fan, J. and Zhang, J. (2000a) Two-step estimation of functional linear models with applications to longitudinal data. *J. R. Statist. Soc. B*, **62**, 303–322.
- Fan, J. and Zhang, W. (1999) Statistical estimation in varying-coefficient models. *Ann. Statist.*, **27**, 1491–1518.
- Fan, J. and Zhang, W. (2000b) Simultaneous confidence bands and hypothesis testing in varying-coefficient models. *Scand. J. Statist.*, **27**, 715–731.

- Gallant, A.R., Hsieh, D.A. and Tauchen, G.E. (1991) On fitting a recalcitrant series: the pound/dollar exchange rate, 1974-1983. In *Nonparametric And Semiparametric Methods in Econometrics and Statistics* (ed. W.A. Barnett, J. Powell and G.E. Tauchen), pp.199-240. Cambridge: Cambridge University Press.
- Härdle, W., Hall, P. and Ichimura, H. (1993) Optimal smoothing in single-index models. *Ann. Statist.*, **21**, 157–178.
- Hastie, T.J. and Tibshirani, R.J. (1993) Varying-coefficient models (with discussion). *J. R. Statist. Soc. B*, **55**, 757–796.
- Hong, Y. and Lee, T.-H. (1999) Inference and forecast of exchange rates via generalized spectrum and nonlinear time series models. *Manuscript*.
- Hoover, D.R., Rice, J.A., Wu, C.O. and Yang, L.-P. (1998) Nonparametric smoothing estimates of time-varying-coefficient models with longitudinal data. *Biometrika*, **85**, 809–822.
- Hurvich, C.M., Simonoff, J.S. and Tsai, C.L. (1998) Smoothing parameter selection in nonparametric regression using an improved Akaike information criterion. *J. Royal Statist. Soc. B*, **60**, 271-293.
- Ichimura, H. (1993) Semiparametric least-squares (SLS) and weighted SLS estimation of single-index models. *J. Econometrics*, **58**, 71–120.
- Kauermann, G. and Tutz, G. (1999) On model diagnostics using varying-coefficient models. *Biometrika*, **86**, 119–128.
- LeBaron, B. (1997) Technical trading rule and regime shifts in foreign exchange. In *Advances in Trading Rules* (ed. E. Acar and S Satchell). Butterworth-Heinemann.
- LeBaron, B. (1999) Technical trading rule profitability and foreign exchange intervention. *J. International Economics*, to appear.
- May, R.M. (1981) Models for two interacting populations. In *Theoretical Ecology* (ed. R.M. May), 78–104. Oxford: Blackwell.
- Newey, W.K. and Stoker, T.M. (1993) Efficiency of weighted average derivative estimators and index models. *Econometrica*, **61**, 1199–1223.
- Nicholls, D.F. and Quinn, B.G. (1982) *Random Coefficient Autoregressive Models: An Introduction*, Lecture Notes in Statistics, No. 11. New York: Springer-Verlag.
- Ramsay, J.O. and Silverman, B.W. (1997) *The Analysis of Functional Data*. New York: Springer-Verlag.
- Ruppert, D. (1997) Empirical-bias bandwidths for local polynomial nonparametric regression and density estimation. *J. Ameri. Statist. Assoc.*, **92**, 1049–1062.
- Samarov, A.M. (1993) Exploring regression structure using nonparametric functional estimation. *J. Ameri. Statist. Assoc.*, **88**, 836–847.
- Seifert, B. and Gasser, Th. (1996) Finite-sample variance of local polynomial: Analysis and solutions. *J. Ameri. Statist. Assoc.*, **91**, 267–275.
- Simonoff, J.S. and Tsai, C.L. (1999) Semiparametric and additive model selection using an improved Akaike information criterion. *Computational and Graphical Statistics*, **8**, 22-40.

- Stenseth, N.C., Falck, W., Bjørnstad, O.N., and Krebs, C.J. (1997) Population regulation in snowshoe hare and Canadian lynx; Asymmetric food web configurations between hare and lynx. *Proceedings of National Academy of Science, US.*, **94**, 5147–5152.
- Stenseth, N.C., Falck, W., Chan, K.S., Bjørnstad, O.N., Tong, H., O'Donoghue, M., Boonstra, R., Boutin, S., Krebs, C.J., and Yoccoz, N.G. (1999) From patterns to processes: phase- and density-dependencies in Canadian lynx cycle. *Proceedings of National Academy of Science, Washington*, to appear.
- Wahba, G. (1977) A survey of some smoothing problems and the method of generalized cross-validation for solving them. In *Applications of Statistics* (ed. P.R. Krishnaiah), 507–523. Amsterdam, North Holland.
- Wu, C.O., Chiang, C.-T. and Hoover, D.R. (1998) Asymptotic confidence regions for kernel smoothing of a varying-coefficient model with longitudinal data. *J. Amer. Statist. Assoc.*, **93**, 1388–1402.
- Xia, Y. and Li, W.K. (1999a) On the estimation and testing of functional-coefficient linear models. *Statistica Sinica*, **9**, 735–757.
- Xia, Y. and Li, W.K. (1999b) On single-index coefficient regression models. *J. Amer. Statist. Assoc.*, **94**, 1275–1285.
- Yao, Q., Tong, H., Finkenstädt, B. and Stenseth, N.C. (2000) Common structure in panels of short time series. *Proc. R. Soc. Lond. B*, **267**, 2459–2467.
- Zhang, W. and Lee, S.Y. (1999) On local polynomial fitting of varying-coefficient models. *Submitted for publication*.
- Zhang, W. and Lee, S.Y. (2000) Variable bandwidth selection in varying-coefficient models. *J. Multivariate Analysis*, to appear.