

# UC Riverside

## UC Riverside Previously Published Works

### Title

Adaptive estimation for varying coefficient models

### Permalink

<https://escholarship.org/uc/item/2pv2z3cv>

### Authors

Chen, Yixin

Wang, Qin

Yao, Weixin

### Publication Date

2015-05-01

### DOI

10.1016/j.jmva.2015.01.017

Peer reviewed

# Adaptive Estimation for Varying Coefficient Models

YIXIN CHEN, QIN WANG AND WEIXIN YAO

## Abstract

In this article, a novel adaptive estimation is proposed for varying coefficient models. Unlike the traditional least squares based methods, the proposed approach can adapt to different error distributions. An efficient EM algorithm is provided to implement the proposed estimation. The asymptotic properties of the resulting estimator are established. Both simulation studies and real data examples are used to illustrate the finite sample performance of the new estimation procedure. The numerical results show that the gain of the new procedure over the least squares estimation can be quite substantial for non-Gaussian errors.

**Key words:** Adaptive estimation; EM algorithm; Kernel smoothing; Local maximum likelihood; Varying coefficient models.

---

<sup>1</sup>Yixin Chen is Senior Specialist Biostatistician, PhD, Department of Biostatistics and Programming, Genzyme Corporation (A Sanofi Company), Cambridge, MA 02142. Email: Yixin.Chen@genzyme.com. Weixin Yao is Associate Professor, Department of Statistics, University of California, Riverside, Riverside, CA 92521. Email: weixin.yao@ucr.edu. Qin Wang is Assistant Professor, Department of Statistical Sciences and Operations Research, Virginia Commonwealth University, Richmond, VA 23284. E-mail: qwang3@vcu.edu.

# 1 Introduction

Since the introduction in Cleveland et al. (1991) and Hastie and Tibshirani (1993), varying coefficient models have gained considerable attention due to their flexibility and good interpretability. They are useful extensions of the classical linear models and have been widely used to explore the dynamic pattern in many scientific areas, such as finance, economics, epidemiology, ecology, etc. By allowing coefficients to vary over the so-called index variable, the modeling bias can be significantly reduced and the ‘curse of dimensionality’ can be avoided (Fan and Zhang, 2008). In recent years, varying coefficient models have experienced rapid developments in both theory and methodology, see, for example, Wu et al. (1998), Hoover et al. (1998), Fan and Zhang (1999, 2000), Cai et al. (2000), Fan and Huang (2005), Wang et al. (2009), Wang and Xia (2009), etc. We refer to readers to Fan and Zhang (2008) for a nice and comprehensive survey.

Let  $y \in \mathcal{R}^1$  be the response,  $\mathbf{x} = (x_1, \dots, x_d)^T \in \mathcal{R}^d$  be the covariate vector, and  $u \in \mathcal{R}^1$  is the index variable. The varying coefficient model is defined as

$$y = \sum_{j=1}^d g_j(u)x_j + \epsilon, \tag{1.1}$$

where  $\{g_1(u), \dots, g_d(u)\}^T$  are unknown smooth coefficient functions. Throughout this article, we assume the random error  $\epsilon$  to be independent of  $(u, \mathbf{x})$ , with mean 0 and a finite second-order moment  $\sigma^2$ . By setting  $x_1 \equiv 1$ , it allows a varying intercept in the model.

Hastie and Tibshirani (1993), Hoover et al. (1998), Chiang et al. (2001) and Eubank et al. (2004) proposed using smoothing spline to estimate coefficient functions. Polynomial spline was used in Huang et al. (2002, 2004) and Huang and Shen (2004). Wu et al. (1998), Hoover et al. (1998), Fan and Zhang (1999) and Kauermann and Tutz (1999) adopted kernel smoothing to estimate coefficient functions. Fan and Zhang

38 (2000) further studied a two-step estimation procedure to deal with the situation where  
39 the coefficient functions admit different degrees of smoothness. Recently, Wang and Xia  
40 (2009) proposed a shrinkage estimation procedure to select important nonparametric  
41 components. Wang et al. (2009) developed a highly robust and efficient procedure based  
42 on local ranks. Nevertheless, most existing methods used least squares type criteria in  
43 estimation, which corresponds to the local likelihood when the error  $\epsilon$  is distributed as  
44 a normal random variable. However, in the absence of normality, the traditional least  
45 squares based estimators will lose some efficiency.

46 In this article, we propose a novel adaptive kernel estimation procedure for varying  
47 coefficient models. It combines the kernel density estimation and the local maximum  
48 likelihood estimation so that the new estimator can adapt to different error distributions.  
49 The new estimator is “adaptive” and “efficient” in the sense that it is asymptotically  
50 equivalent to the infeasible local likelihood estimator (Staniswalis, 1989; Fan et al., 1998),  
51 which requires the knowledge of the error distribution. An efficient EM algorithm is  
52 proposed to implement the adaptive estimation. We demonstrate through a simulation  
53 study that the new estimate is more efficient than the existing least squares based  
54 kernel estimate when the error distribution deviates from normal. In addition, when the  
55 error is exactly normal, the new method is broadly comparable to the existing kernel  
56 approach. We further illustrate the effectiveness of the proposed adaptive estimation  
57 method through two real data examples.

58 The rest of the article is organized as follows. In section 2, we introduce the new  
59 adaptive estimation for the varying coefficient models and the EM algorithm. In section  
60 3, we compare our proposed approach with the traditional least squares based estimation  
61 for five different error densities through a simulation study and then apply the new  
62 method to two real data examples. We conclude this article with a brief discussion in  
63 Section 4. All technical conditions and proofs are relegated to the Appendix.

## 2 New Adaptive Estimation

### 2.1 Introduction to the new method

Suppose that  $\{\mathbf{x}_i, u_i, y_i, i = 1, \dots, n\}$  is a random sample from model (1.1). For  $u$  in a neighborhood of  $u_0$ , we can approximate the varying coefficient functions locally as

$$g_j(u) \approx g_j(u_0) + g'_j(u_0)(u - u_0) \equiv b_j + c_j(u - u_0), \quad \text{for } j = 1, \dots, d. \quad (2.1)$$

The traditional local linear estimation of (1.1) is to minimize

$$\sum_{i=1}^n K_h(u_i - u_0) \left[ y_i - \sum_{j=1}^d \{b_j + c_j(u_i - u_0)\} x_{ij} \right]^2, \quad (2.2)$$

with respect to  $(b_1, \dots, b_d)$  and  $(c_1, \dots, c_d)$  for a given kernel density  $K(\cdot)$  and a bandwidth  $h$ , where  $K_h(t) = h^{-1}K(t/h)$ . It is well known that the choice of kernel function is not critical in terms of estimation efficiency. Throughout this article, a Gaussian kernel will be used for  $K(\cdot)$ . Due to the least squares in (2.2), the resulting estimate may lose some efficiency when the error distribution is not normal. Therefore, it is desirable to develop an estimation procedure which can adapt to different error distributions.

Let  $f(\epsilon)$  be the density function of  $\epsilon$ . If  $f(\epsilon)$  were known, it would be natural to estimate the local parameters in (2.1) by maximizing the following local log-likelihood function

$$\sum_{i=1}^n K_h(u_i - u_0) \log f \left[ y_i - \sum_{j=1}^d \{b_j + c_j(u_i - u_0)\} x_{ij} \right]. \quad (2.3)$$

However, in practice,  $f(\epsilon)$  is generally unknown but can be replaced by a kernel density estimate based on the initial estimated residual  $\tilde{\epsilon}_1, \dots, \tilde{\epsilon}_n$ ,

$$\tilde{f}(\epsilon_i) = \frac{1}{n} \sum_{j \neq i}^n K_{h_0}(\epsilon_i - \tilde{\epsilon}_j), \quad \text{for } i, j = 1, 2, \dots, n \quad (2.4)$$

80 where  $\tilde{\epsilon}_i = y_i - \sum_{j=1}^d \tilde{g}_j(u_i)x_{ij}$  and  $\tilde{g}_j(\cdot)$  can be estimated by least squares (or  $L_1$  norm,  
81 i.e., median regression) based local linear estimate (2.2). Here we use leave-one-out kernel  
82 density estimate for  $f(\epsilon_i)$  to remove the estimation bias. Let  $\boldsymbol{\theta} = (b_1, \dots, b_d, c_1, \dots, c_d)^T$ .  
83 Then our proposed adaptive local linear estimate for the local parameter  $\boldsymbol{\theta}$  is

$$\hat{\boldsymbol{\theta}} = \arg \max_{\boldsymbol{\theta}} Q(\boldsymbol{\theta}), \quad (2.5)$$

84 where

$$Q(\boldsymbol{\theta}) = \sum_{i=1}^n K_h(u_i - u_0) \log \left( \frac{1}{n} \sum_{j \neq i} K_{h_0} \left[ y_i - \sum_{l=1}^d \{b_l + c_l(u_i - u_0)\} x_{il} - \tilde{\epsilon}_j \right] \right). \quad (2.6)$$

85 The idea of adaptiveness can be traced back to Beran (1974) and Stone (1975),  
86 where the adaptive estimation was proposed for location models. Later, this idea was  
87 extended to regression, time series and other models, see Bickel (1982), Manski (1984),  
88 Steigerwald (1992), Schick (1993), Drost and Klaassen (1997), Hodgson (1998), Yuan  
89 and De Gooijer (2007), and Yuan (2009). Linton and Xiao (2007) proposed an elegant  
90 adaptive nonparametric regression estimator by maximizing the local likelihood function.  
91 In fact, the adaptive method proposed in Linton and Xiao (2007) can be seen as a special  
92 case of ours when  $d = 1$  in (1.1). Recently, Wang and Yao (2012) and Yao and Zhao  
93 (2013) extended the idea of adaptive estimation to sufficient dimension reduction and  
94 linear regression, respectively.

## 95 **2.2 Computation: an EM algorithm**

96 Unlike least squares criterion, (2.5) does not have an explicit solution due to the sum-  
97 mation inside the log function, which is similar to the mixture structure. In this section,  
98 we propose an EM algorithm to maximize it by extending the generalized modal EM  
99 algorithm proposed in Yao (2013).

100 Let  $\boldsymbol{\theta}^{(0)}$  be an initial parameter estimate, such as the least squares (or  $L_1$  norm, i.e.,  
 101 median regression) based local linear estimate. We can update the parameter estimate  
 102 according to the following algorithm.

103 **Algorithm 2.1.** *At  $(k + 1)$ th step, we calculate the following E and M steps:*

E-Step: Calculate the classification probabilities  $p_{ij}^{(k+1)}$ ,

$$\begin{aligned}
 p_{ij}^{(k+1)} &= \frac{K_{h_0} \left[ y_i - \sum_{l=1}^d \{b_l^{(k)} + c_l^{(k)}(u_i - u_0)\} x_{il} - \tilde{\epsilon}_j \right]}{\sum_{j \neq i} K_{h_0} \left[ y_i - \sum_{l=1}^d \{b_l^{(k)} + c_l^{(k)}(u_i - u_0)\} x_{il} - \tilde{\epsilon}_j \right]} \\
 &\propto K_{h_0} \left[ y_i - \sum_{l=1}^d \{b_l^{(k)} + c_l^{(k)}(u_i - u_0)\} x_{il} - \tilde{\epsilon}_j \right], \quad 1 \leq j \neq i \leq n. \quad (2.7)
 \end{aligned}$$

M-Step: Update  $\boldsymbol{\theta}^{(k+1)}$ ,

$$\begin{aligned}
 \boldsymbol{\theta}^{(k+1)} &= \arg \max_{\boldsymbol{\theta}} \sum_{i=1}^n \sum_{j \neq i} \left\{ p_{ij}^{(k+1)} K_h(u_i - u_0) \log \left( K_{h_0} \left[ y_i - \sum_{l=1}^d \{b_l + c_l(u_i - u_0)\} x_{il} - \tilde{\epsilon}_j \right] \right) \right\} \\
 &= \arg \min_{\boldsymbol{\theta}} \sum_{i=1}^n \sum_{j \neq i} \left\{ p_{ij}^{(k+1)} K_h(u_i - u_0) [y_i - \tilde{\epsilon}_j - \mathbf{z}_i^T \boldsymbol{\theta}]^2 \right\}, \\
 &= \left( \sum_{i=1}^n \sum_{j \neq i} p_{ij}^{(k+1)} K_h(u_i - u_0) \mathbf{z}_i \mathbf{z}_i^T \right)^{-1} \sum_{i=1}^n \sum_{j \neq i} p_{ij}^{(k+1)} K_h(u_i - u_0) (y_i - \tilde{\epsilon}_j) \mathbf{z}_i
 \end{aligned} \quad (2.8)$$

104 where  $\mathbf{z}_i = \{\mathbf{x}_i^T, \mathbf{x}_i^T(u_i - u_0)\}^T$  and the second equation follows the use of Gaussian  
 105 kernel for density estimation.

106 The above EM algorithm monotonically increases the estimated local log-likelihood  
 107 (2.6) after each iteration, as shown in the following proposition. Its proof is given in the  
 108 appendix.

**Proposition 2.1.** *Each iteration of the above E and M steps will monotonically*

increase the local log-likelihood (2.6), i.e.,

$$Q(\boldsymbol{\theta}^{(k+1)}) \geq Q(\boldsymbol{\theta}^{(k)}),$$

109 for all  $k$ , where  $Q(\cdot)$  is defined as in (2.6).

## 110 2.3 Asymptotic result

We now establish the consistency and derive the asymptotic distribution of the proposed adaptive local linear estimator of  $\boldsymbol{\theta}$ . Define  $\mu_k = \int u^k K(u) du$  and  $\nu_k = \int u^k K^2(u) du$ . Let  $\mathbf{H} = \text{diag}(1, h) \otimes \mathbf{I}_d$  with  $\otimes$  denoting the Kronecker product and  $\mathbf{I}_d$  being the  $d \times d$  identity matrix. Let  $q(\cdot)$  denote the marginal density of  $u$ , and

$$\Gamma_{jk}(u_i) = \mathbb{E}(x_{ij}x_{ik}|u_i) \text{ for } 1 \leq j, k \leq d, i = 1, \dots, n, \quad (2.9)$$

$$\boldsymbol{\Gamma}(u_0) = \{\Gamma_{jk}(u_0)\}_{1 \leq j, k \leq d}. \quad (2.10)$$

**Theorem 2.1.** *Under the regularity conditions in the Appendix, with probability approaching 1, there exists a consistent local maximizer  $\hat{\boldsymbol{\theta}} = (\hat{b}_1, \dots, \hat{b}_d, \hat{c}_1, \dots, \hat{c}_d)^T$  of (2.6) such that*

$$\mathbf{H}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}) = O_p\{(nh)^{-1/2} + h^2\}.$$

111

112 Based on Theorem 2.1, we can know that the proposed adaptive local linear estimator  
113 of  $\boldsymbol{\theta}$  is consistent and its proof is provided in the Appendix. Next, we provide the  
114 asymptotic distribution of the proposed estimator.

**Theorem 2.2.** *Suppose that the regularity conditions in the Appendix hold. Then*



$\hat{\boldsymbol{\theta}}$ , given in Theorem 2.1, has the following asymptotic distribution

$$\sqrt{nh} \left\{ \mathbf{H}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}) - \mathbf{S}^{-1} \frac{h^2}{2} \sum_{j=1}^d g_j''(u_0) \boldsymbol{\psi}_j (1 + o_p(1)) \right\} \xrightarrow{D} N(\mathbf{0}_{2d}, [E\{\rho'(\epsilon_i)^2\}]^{-1} q(u_0)^{-1} \mathbf{S}^{-1} \boldsymbol{\Lambda} \mathbf{S}^{-1}),$$

115 where  $\mathbf{0}_{2d}$  is a  $2d \times 1$  vector with each entry being 0,  $\rho(\cdot) = \log f(\cdot)$ ,  $\mathbf{S} = \begin{pmatrix} 1 & 0 \\ 0 & \mu_2 \end{pmatrix} \otimes$

116  $\boldsymbol{\Gamma}(u_0)$ ,  $\boldsymbol{\Lambda} = \begin{pmatrix} \nu_0 & \nu_1 \\ \nu_1 & \nu_2 \end{pmatrix} \otimes \boldsymbol{\Gamma}(u_0)$ ,  $\boldsymbol{\psi}_j = \begin{pmatrix} \mu_2 \\ \mu_3 \end{pmatrix} \otimes (\boldsymbol{\Gamma}_{jk}(u_0))_{1 \leq k \leq d}^T$ , and  $\boldsymbol{\Gamma}(u_0)$  is given

117 by (2.10).

118 A sketch of the proof of the above theorems is provided in the Appendix. As shown  
119 in Linton and Xiao (2007), one important property of the proposed adaptive estimate  
120 is that it achieves the same asymptotic efficiency as if the error density were known.  
121 Therefore, estimating  $f$  by kernel density estimation will not affect the asymptotic dis-  
122 tribution of the resulting estimator of  $\boldsymbol{\theta}$ . As Linton and Xiao (2007) pointed out that  
123 such a new estimation method can “do as well as the corresponding estimator one would  
124 compute if one knew the error density.” However it is not possible to achieve the lower  
125 bound here (Fan, 1993). Any specific estimator can be bettered for some specific model  
126 setting.

127 Note that the least squares based local linear estimate (Zhang and Lee, 2000), by  
128 minimizing (2.2), has the same asymptotic bias as the new method but slightly different  
129 asymptotic variance, which replaces  $[E\{\rho'(\epsilon_i)^2\}]^{-1}$  by  $\sigma^2 = E(\epsilon^2)$ . Based on Cauchy-  
130 Schwarz inequality, we have

$$E\{\rho'(\epsilon_i)^2\}E(\epsilon^2) \geq [E\{\epsilon\rho'(\epsilon)\}]^2 = 1$$

131 and the equality holds if and only if  $\rho'(\epsilon) \propto \epsilon$ , i.e.,  $f(\epsilon)$  is a normal density. Therefore,

132  $[E\{\rho'(\epsilon_i)^2\}]^{-1} \leq \sigma^2$  and the asymptotic variance of the new estimator is no larger than  
 133 that of least squares based local linear estimate for any error density  $f(\epsilon)$ .

## 134 3 Examples

### 135 3.1 Simulation study

136 In this section, we conduct a simulation study to compare the proposed adaptive es-  
 137 timation (Adapt) with the traditional least squares based kernel estimation (LS) for  
 138 varying coefficient models. The following five error distributions of  $\epsilon$  were considered in  
 139 our numerical experiment:

- 140 1.  $N(0, 1)$ ;
- 141 2.  $t_3$ ;
- 142 3.  $0.5N(-1, 0.5^2) + 0.5N(1, 0.5^2)$ ;
- 143 4.  $0.3N(-1.4, 1) + 0.7N(0.6, 0.4^2)$ ;
- 144 5.  $0.9N(0, 1) + 0.1N(0, 10^2)$ .

145 The standard normal distribution serves as a baseline in our comparison. The second  
 146 one is a  $t$ -distribution with 3 degrees of freedom. The third density is bimodal and the  
 147 fourth one is left skewed. The last one is a contaminated normal mixture distribution,  
 148 where 10% of the data from  $N(0, 10^2)$  are most likely to be outliers.

149 For each of the above error distributions, we consider the following two models:

150 **Model 1:**  $y = g_1(u) + g_2(u)x_1 + g_3(u)x_2 + \epsilon$ , where  $g_1(u) = \exp(2u - 1)$ ,  $g_2(u) = 8u(1 - u)$ ,  
 151 and  $g_3(u) = 2 \sin^2(2\pi u)$ .

152 **Model 2:**  $y = g_1(u) + g_2(u)x_1 + g_3(u)x_2 + \epsilon$ , where  $g_1(u) = \sin(2\pi u)$ ,  $g_2(u) = (2u -$   
 153  $1)^2 + 0.5$ , and  $g_3(u) = \exp(2u - 1) - 1$ .

154 In both models,  $x_1$  and  $x_2$  follow a standard normal distribution with correlation co-  
155 efficient  $\gamma = 1/\sqrt{2}$ . The index variable  $u$  is a uniform random variable on  $[0, 1]$ , and  
156 is independent of  $(x_1, x_2)$ . There are two bandwidths in the estimation,  $h$  in the local  
157 log-likelihood and  $h_0$  in the kernel density estimation. An asymptotic optimal  $h$  can  
158 be found by minimizing the asymptotic mean squared errors provide in Theorem 2.2  
159 and can be estimated by a plug-in estimator which replaces the unknown quantities in  
160 Theorem 2.2 by their estimates. In our examples, the bandwidth  $h$  is chosen by leave-  
161 one-out cross-validation with more details in Fan and Zhang (1999), and  $h_0 = h/\log(n)$   
162 following Linton and Xiao (2007). The performance of estimator  $\hat{g}(\cdot)$  is assessed via the  
163 square root of the average squared errors (RASE; Cai et al., 2000; Wang et al., 2009),

$$\text{RASE}^2 = \frac{1}{N} \sum_{k=1}^N \sum_{j=1}^3 [\hat{g}_j(u_k) - g_j(u_k)]^2, \quad (3.1)$$

164 where  $u_k$ ,  $k = 1, \dots, N$ , are the equally spaced grid points at which the functions  $g_j(\cdot)$   
165 were evaluated. We conduct two sets of simulations with sample size  $n=200$  and 400  
166 respectively, each with 200 data replications.

167 The simulation results are summarized in Tables 1 and 2. We can clearly see that the  
168 proposed adaptive estimation outperforms the least squares method when the error is  
169 non-normal. The gain in estimation efficiency can be quite substantial even for moderate  
170 sample sizes. When the error follows exactly normal distribution, our approach is still  
171 broadly comparable with the least squares based method.

172 Figures 1 and 2 plot the estimated coefficient functions and the 95% pointwise confi-  
173 dence intervals based on a typical sample when  $n=200$  and the error distribution is the  
174 contaminated normal mixture (Case 5). Due to the complex forms of the asymptotic  
175 standard errors of the coefficient functions, similar to Wang, Kai and Li (2009), we adopt  
176 the bootstrap method to calculate the 95% pointwise confidence intervals. As expected,

177 the adaptive estimation method provides narrower confidence intervals than the least  
 178 squares based method, since the adaptive method provides more accurate estimate than  
 179 the least squares estimate when the error is not normal.

Table 1: Model 1 estimation accuracy comparison–RASE and its standard error in brackets.

$\epsilon$	$n = 200$		$n = 400$	
	LS	Adapt	LS	Adapt
1	0.483(0.079)	0.439(0.081)	0.366(0.053)	0.324(0.053)
2	0.671(0.167)	0.601(0.139)	0.493(0.111)	0.422(0.086)
3	0.500(0.083)	0.401(0.077)	0.379(0.061)	0.277(0.048)
4	0.508(0.088)	0.376(0.082)	0.383(0.062)	0.262(0.045)
5	1.188(0.411)	0.720(0.220)	0.871(0.227)	0.459(0.098)

Table 2: Model 2 estimation accuracy comparison–RASE and its standard error in brackets.

$\epsilon$	$n = 200$		$n = 400$	
	LS	Adapt	LS	Adapt
1	0.362(0.077)	0.380(0.074)	0.263(0.051)	0.275(0.049)
2	0.618(0.301)	0.566(0.201)	0.431(0.129)	0.384(0.076)
3	0.412(0.091)	0.351(0.080)	0.290(0.059)	0.215(0.041)
4	0.407(0.102)	0.319(0.089)	0.291(0.061)	0.207(0.051)
5	1.133(0.397)	0.669(0.224)	0.828(0.224)	0.436(0.101)

### 180 3.2 Real-data applications

*Example 1 (Hong Kong environmental data).* We now illustrate the adaptive estimation method via an application to an environmental data set. The data were collected daily in Hong Kong from January 1, 1994, to December 31, 1995 and have been analyzed by

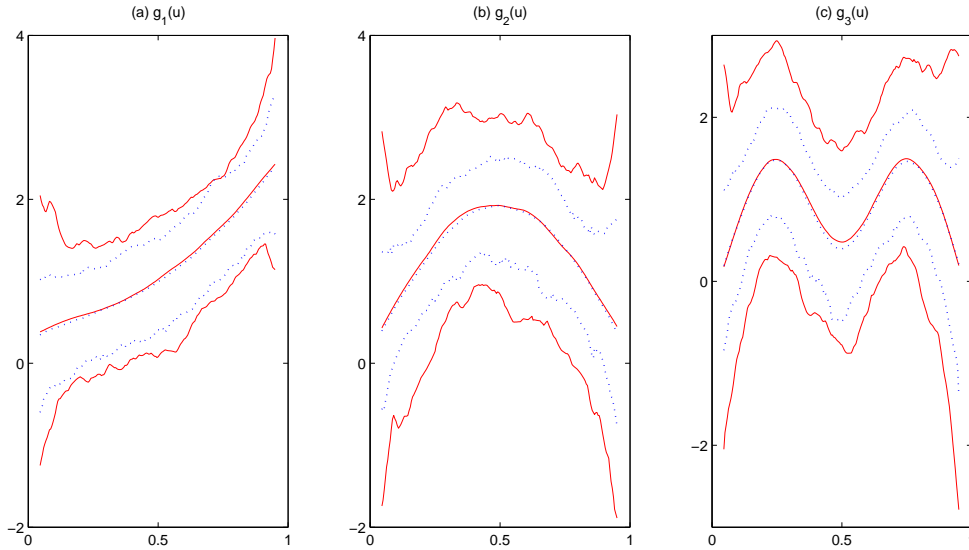


Figure 1: Estimated coefficient functions with 95% pointwise confidence intervals (blue dotted line for Adapt and red solid line for LS) for model 1.

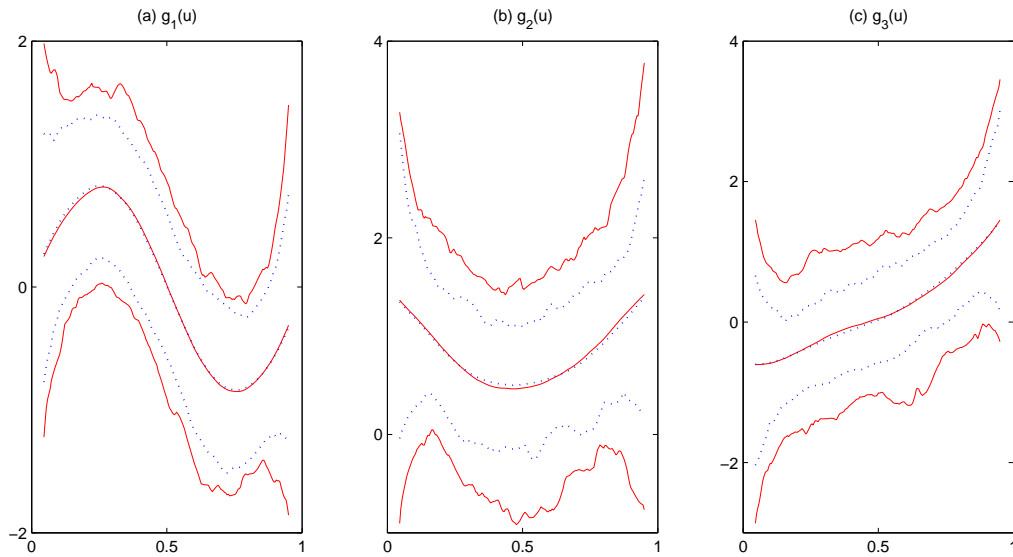


Figure 2: Estimated coefficient functions with 95% pointwise confidence intervals (blue dotted line for Adapt and red solid line for LS) for model 2.

Fan and Zhang (1999), Cai et al. (2000), Xia et al. (2002) and Fan and Zhang (2008). In this data set, a collection of daily measurements of pollutants and other environmental factors are included. Following Fan and Zhang (1999), we consider three pollutants:

sulphur dioxide  $x_2$  (in  $\mu g/m^3$ ), nitrogen dioxide  $x_3$  (in  $\mu g/m^3$ ), and respirable suspended particulates  $x_4$  (in  $\mu g/m^3$ ) (this variable is named as ‘dust’ in Fan and Zhang (1999), Fan and Zhang (2008), and Cai et al. (2000)). The response variable  $y$  is the logarithm of the number of daily hospital admissions. We set  $x_1 = 1$  as the intercept term and let  $u$  denote time which is scaled to the interval  $[0, 1]$ . As in the previous analyses, all three predictors are centered. The following varying coefficient model is considered to investigate the relationship between  $y$  and the levels of pollutants  $x_2$ ,  $x_3$ , and  $x_4$ .

$$y = g_1(u) + g_2(u)x_2 + g_3(u)x_3 + g_4(u)x_4 + \epsilon.$$

181 We set aside 50 observations as the test set. The bandwidth  $h$ , selected by leave-  
 182 one-out cross-validation, is around 0.146. The estimated coefficient functions together  
 183 with 95% pointwise confidence intervals are depicted in Figure 3. We also compare  
 184 the median squared prediction errors,  $MSPE = \text{Median}\{(y_j - \hat{y}_j)^2, j = 1, \dots, k\}$ , from  
 185 our adaptive approach and the traditional least squares estimation, where  $k = 50$  and  
 186  $\hat{y}_j = \hat{g}_1(u_j) + \hat{g}_2(u_j)x_{j2} + \hat{g}_3(u_j)x_{j3} + \hat{g}_4(u_j)x_{j4}$ . The MSPE from our adaptive approach  
 187 is 0.0183, compared to 0.0178 from the LS estimation.

188 In Figure 5 (a), we give the residual QQ-plot for Hong Kong environmental data.  
 189 From the plot, we can see that the residual is very close to normal, which explains why  
 190 the MSPE of the adaptive approach is close to the MSPE of the LS estimation.

*Example 2 (Boston housing data).* The Boston Housing Data (corrected version in Gilley and Pace (1996)), which has been analyzed by Fan and Huang (2005), Wang and Xia (2009) and Sun et al. (2014), is publicly available in the R package *mlbench*, (<http://cran.r-project.org/>). This data set includes the median value of owner-occupied homes in 506 U.S. census tracts of the Boston area in 1970 and several variables that might explain the variation of housing values. Following Fan and Huang (2005) and Wang and Xia (2009), we considered seven independent variables: CRIM (per capita

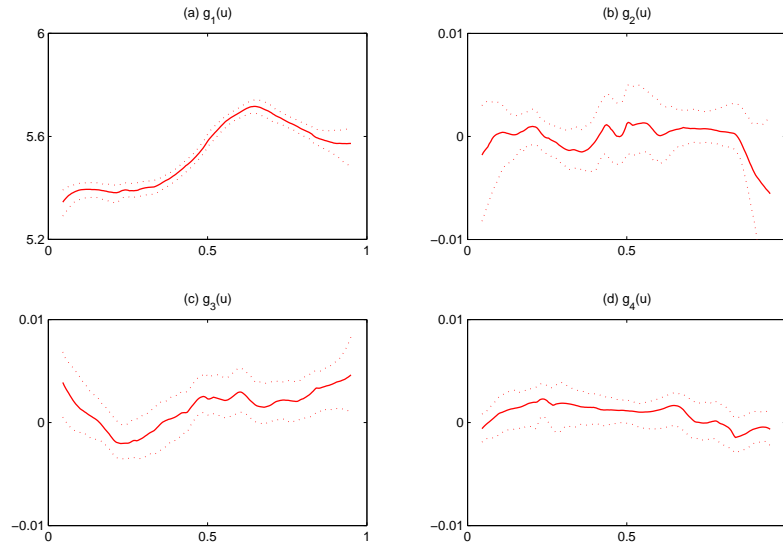


Figure 3: Estimated coefficient functions (solid curves) with 95% pointwise confidence intervals (dotted curves) for Hong Kong environmental data.

crime rate by town), RM (average number of rooms per dwelling), TAX (full-value property-tax rate per \$10,000), NOX (nitric oxides concentration parts per 10 million), PTRATIO (pupil-teacher ratio by town), AGE (proportion of owner-occupied units built prior to 1940), and LSTAT (lower status of the population). The response variable is CMEDV (corrected median value of owner-occupied homes in USD 1000's). We denote the covariates CRIM, RM, TAX, NOX, PTRATIO and AGE to be  $x_2, x_3, \dots, x_7$ , respectively. Let  $x_1 = 1$  be the intercept term and  $u = \sqrt{\text{LSTAT}}$  be the index variable. By doing so, we can fit different regression models at different lower status population percentage (Fan and Huang, 2005). Following Fan and Huang (2005) we use the square root transformation on the index variable LSTAT to make the data symmetrically distributed. The following varying coefficient model was fit to the data,

$$y_i = g_1(u_i) + \sum_{j=2}^7 g_j(u_i)x_{ij} + \epsilon_i.$$

191 Similar to the analysis in the previous example, we set aside 50 observations for check-  
 192 ing prediction errors. The bandwidth  $h$  was selected by leave-one-out cross-validation,  
 193 which is around 0.294. The estimated coefficient functions are depicted in Figure 4.  
 194 From the plot, we can see that the coefficient functions of  $x_2$  (CRIM) and  $x_3$  (RM) vary  
 195 over time. The coefficient functions of  $x_4$  (TAX),  $x_5$  (NOX), and  $x_7$  (AGE) are very  
 196 close to zero and the coefficient function of  $x_6$  (PTRATIO) shows no significant trend.  
 197 These discoveries are consistent with those from Fan and Huang (2005) and Wang and  
 198 Xia (2009). In terms of the median squared prediction error (MSPE), the MSPE from  
 199 our adaptive approach is 0.0484, compared to 0.0604 from the LS estimation.

200 In Figure 5 (b), the QQ-plot of residuals from the above fit showed a clear deviation  
 201 from normality, which explains why the MSPE from the adaptive approach is much  
 202 smaller than the MSPE from the LS estimation.

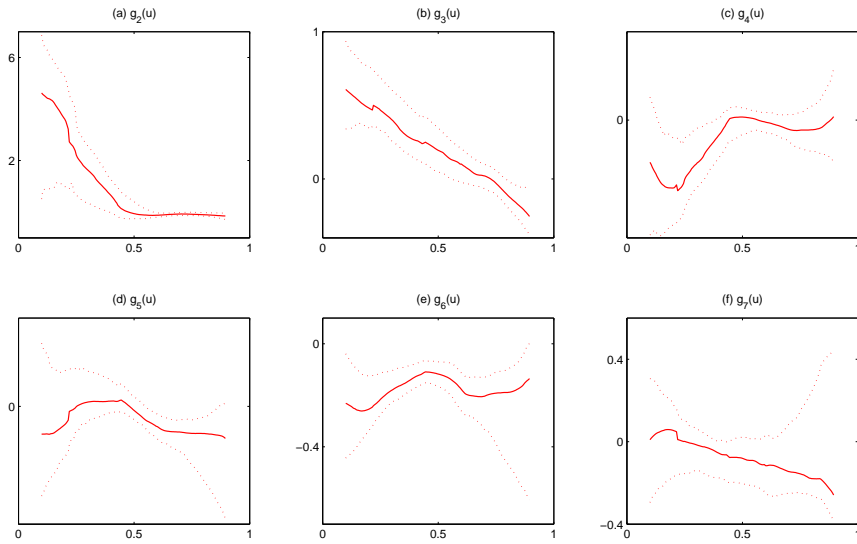


Figure 4: Estimated coefficient functions (solid curves) with 95% pointwise confidence intervals (dotted curves) for Boston housing data.



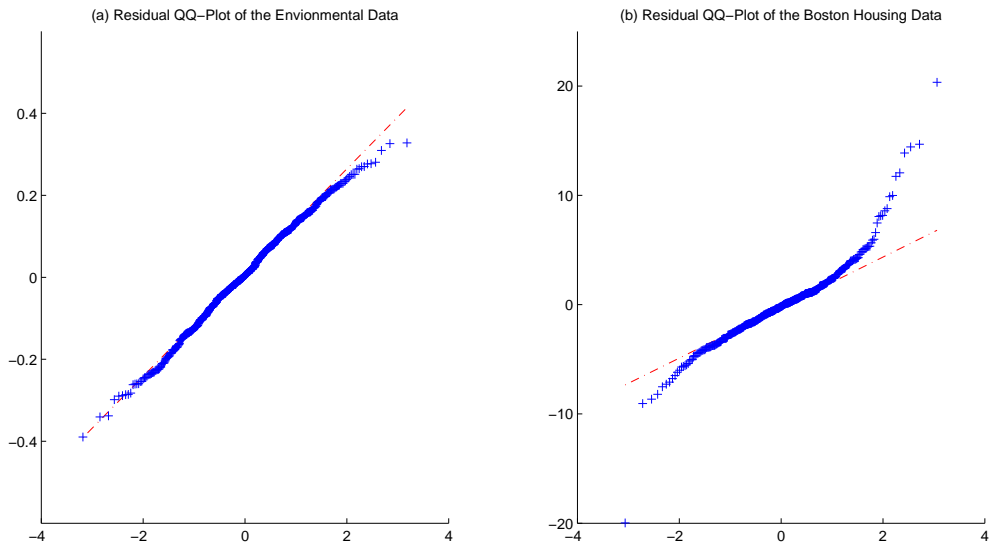


Figure 5: Residual QQ-plot for two data examples: (a) Hong Kong environmental data; (b) Boston housing data.

## 203 4 Discussion

204 In this article, we proposed an adaptive estimation for varying coefficient models. The  
 205 new estimation procedure can adapt to different errors and thus provide a more efficient  
 206 estimate than the traditional least squares based estimate. Simulation studies and two  
 207 real data applications confirmed our theoretical findings.

208 It will be interesting to know whether we can also perform some adaptive hypothesis  
 209 tests for the coefficient functions using the estimated error density. For example, we  
 210 might be interested in testing some parametric assumptions, such as constant or zero, for  
 211 the coefficient functions. It requires more research about whether the Wilks phenomenon  
 212 for generalized likelihood ratio statistic proposed by Fan et al. (2001) still holds for the  
 213 proposed adaptive varying coefficient models.

214 The idea of the proposed adaptive estimator might also be generalized to many other  
 215 models, such as varying coefficient partial linear models and nonparametric additive  
 216 models. In addition, by combining this adaptive idea with shrinkage estimation, we can

217 develop adaptive variable selection procedures. Such study is under way.

218 Zhang and Lee (2000) investigated variable bandwidth selection for varying coef-  
219 ficient models and studied asymptotic properties of the resulting estimators and the  
220 bandwidth. It is our interest to extend their variable bandwidth selection method and  
221 the corresponding asymptotic properties to our adaptive estimation procedure.

As one referee pointed out that we could also extend the idea of Yuan and De  
Gooijer (2007) to derive an adaptive estimate for varying coefficient model. Let  $\epsilon_i(\boldsymbol{\theta}) =$   
 $y_i - \sum_{l=1}^d [b_l + c_l(u_i - u_0)]$ , and

$$f_n(\epsilon_i(\boldsymbol{\theta})) = \frac{1}{n-1} \sum_{j \neq i} K_h(r(\epsilon_i(\boldsymbol{\theta})) - r(\epsilon_j(\boldsymbol{\theta}))).$$

Based on Yuan and De Gooijer (2007), we can estimate  $\boldsymbol{\theta}$  by

$$\hat{\boldsymbol{\theta}} = \arg \max_{\boldsymbol{\theta}} \sum_{i=1}^n K_h(u_i - u_0) \log f_n(\epsilon_i(\boldsymbol{\theta})).$$

222 Here,  $r(\cdot)$  is some monotone nonlinear function that is used to avoid the cancelation  
223 of the intercept terms  $b_l$ s in  $f_n(\epsilon_i(\boldsymbol{\theta}))$ . One advantage of the above method is that it  
224 does not require an initial estimate. However, compared to the proposed estimate in  
225 this paper, the asymptotic variance of the above estimator depends on the choice of  $r(\cdot)$   
226 and generally does not reach the Cramér-Rao lower bound for a nonlinear function  $r(\cdot)$ .  
227 In addition, the computation of the above estimator is also more expensive due to the  
228 nonlinear function  $r(\cdot)$ .

## 229 Acknowledgements

230 The authors thank the editor, the associate editor, and reviewers for their constructive  
231 comments that have led to a dramatic improvement of the earlier version of this article.

232 Yao's research is supported by NSF grant DMS-1461677.

## 233 Appendix

234 We first list the regularity conditions used in our proof.

### 235 Conditions:

- 236 1.  $K(\cdot)$  is bounded, symmetric, and has bounded support and bounded derivative;
- 237 2.  $\{\mathbf{x}_i\}_i$ ,  $\{u_i\}_i$ ,  $\{\epsilon_i\}_i$  are independent and identically distributed and  $\{\epsilon_i\}_i$  is independent of  $\{\mathbf{x}_i\}_i$  and  $\{u_i\}_i$ , where  $\{\mathbf{x}_i\}_i$  means  $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ , same for notations  $\{u_i\}_i$  and  $\{\epsilon_i\}_i$ . Additionally, the predictor  $\mathbf{x}$  has a bounded support;
- 240 3. The probability distribution function  $f(\cdot)$  of  $\epsilon$  has bounded continuous derivatives up to order 4. Let  $\rho(\epsilon) = \log f(\epsilon)$ . Assume  $E[\rho'(\epsilon_i)] = 0$ ,  $E[\rho''(\epsilon_i)] < \infty$ ,  
241  $E[\rho'(\epsilon_i)^2] < \infty$  and  $\rho'''(\cdot)$  is bounded;
- 243 4. The marginal density of  $u$  has a continuous second derivative in some neighborhood  
244 of  $u_0$  and  $q(u_0) \neq 0$ ;
- 245 5.  $h \rightarrow 0$ ,  $nh \rightarrow \infty$  as  $n \rightarrow \infty$  and  $h_0 = h/\log(n)$ ;
- 246 6.  $g_j(\cdot)$  has bounded, continuous  $3^{rd}$  derivatives for  $1 \leq j \leq d$ .

247 These conditions are adopted from Fan and Zhang (1999) and Linton and Xiao (2007).  
248 They are not the weakest possible conditions. For instance, we can relax the bounded  
249 support assumption of  $K(\cdot)$ . All the asymptotic results still hold if we put a restriction  
250 on the tail of  $K(\cdot)$ . For example,  $\limsup_{t \rightarrow \infty} |K(t)t^5| < \infty$  (Fan and Gijbels, 1992).  
251 The independence of  $\{\mathbf{x}_i\}_i$  and  $\{\epsilon_i\}_i$  can be relaxed based on the discussion of Section  
252 4 of Linton and Xiao (2007).

Note that

$$\begin{aligned}
 & Q(\boldsymbol{\theta}^{(k+1)}) - Q(\boldsymbol{\theta}^{(k)}) \\
 &= \sum_{i=1}^n K_h(u_i - u_0) \log \left\{ \frac{\sum_{j \neq i} K_{h_0} \left[ y_i - \sum_{l=1}^d \left\{ b_l^{(k+1)} + c_l^{(k+1)}(u_i - u_0) \right\} x_{il} - \tilde{\epsilon}_j \right]}{\sum_{j \neq i} K_{h_0} \left[ y_i - \sum_{l=1}^d \left\{ b_l^{(k)} + c_l^{(k)}(u_i - u_0) \right\} x_{il} - \tilde{\epsilon}_j \right]} \right\} \\
 &= \sum_{i=1}^n K_h(u_i - u_0) \log \sum_{j \neq i} \left( \frac{K_{h_0} \left[ y_i - \sum_{l=1}^d \left\{ b_l^{(k)} + c_l^{(k)}(u_i - u_0) \right\} x_{il} - \tilde{\epsilon}_j \right]}{\sum_{j \neq i} K_{h_0} \left[ y_i - \sum_{l=1}^d \left\{ b_l^{(k)} + c_l^{(k)}(u_i - u_0) \right\} x_{il} - \tilde{\epsilon}_j \right]} \right) \\
 &\quad \times \left( \frac{K_{h_0} \left[ y_i - \sum_{l=1}^d \left\{ b_l^{(k+1)} + c_l^{(k+1)}(u_i - u_0) \right\} x_{il} - \tilde{\epsilon}_j \right]}{K_{h_0} \left[ y_i - \sum_{l=1}^d \left\{ b_l^{(k)} + c_l^{(k)}(u_i - u_0) \right\} x_{il} - \tilde{\epsilon}_j \right]} \right) \\
 &= \sum_{i=1}^n K_h(u_i - u_0) \log \left\{ \sum_{j \neq i} p_{ij}^{(k+1)} \frac{K_{h_0} \left[ y_i - \sum_{l=1}^d \left\{ b_l^{(k+1)} + c_l^{(k+1)}(u_i - u_0) \right\} x_{il} - \tilde{\epsilon}_j \right]}{K_{h_0} \left[ y_i - \sum_{l=1}^d \left\{ b_l^{(k)} + c_l^{(k)}(u_i - u_0) \right\} x_{il} - \tilde{\epsilon}_j \right]} \right\},
 \end{aligned}$$

where

$$p_{ij}^{(k+1)} = \frac{K_{h_0} \left[ y_i - \sum_{l=1}^d \left\{ b_l^{(k)} + c_l^{(k)}(u_i - u_0) \right\} x_{il} - \tilde{\epsilon}_j \right]}{\sum_{j \neq i} K_{h_0} \left[ y_i - \sum_{l=1}^d \left\{ b_l^{(k)} + c_l^{(k)}(u_i - u_0) \right\} x_{il} - \tilde{\epsilon}_j \right]}.$$

From the Jensen's inequality, we have

$$\begin{aligned}
 & Q(\boldsymbol{\theta}^{(k+1)}) - Q(\boldsymbol{\theta}^{(k)}) \\
 &\geq \sum_{i=1}^n K_h(u_i - u_0) \sum_{j \neq i} p_{ij}^{(k+1)} \log \left\{ \frac{K_{h_0} \left[ y_i - \sum_{l=1}^d \left\{ b_l^{(k+1)} + c_l^{(k+1)}(u_i - u_0) \right\} x_{il} - \tilde{\epsilon}_j \right]}{K_{h_0} \left[ y_i - \sum_{l=1}^d \left\{ b_l^{(k)} + c_l^{(k)}(u_i - u_0) \right\} x_{il} - \tilde{\epsilon}_j \right]} \right\}.
 \end{aligned}$$

254 Based on the property of M-step of (2.8), we have  $Q(\boldsymbol{\theta}^{(k+1)}) - Q(\boldsymbol{\theta}^{(k)}) \geq 0$ .  $\square$

255 **Proof of Theorem 2.1**

256 Note that the estimator  $\hat{\boldsymbol{\theta}}$  is the maximizer of the following objective function

$$\arg \max_{\boldsymbol{\theta}} \sum_{i=1}^n K_h(u_i - u_0) \log \tilde{f} \left[ y_i - \sum_{l=1}^d \{b_l + c_l(u_i - u_0)\} x_{il} \right], \quad (4.1)$$

where

$$\tilde{f}(\epsilon_i) = \frac{1}{n} \sum_{j \neq i} K_{h_0}(\epsilon_i - \tilde{\epsilon}_j)$$

257 is the kernel density estimate of  $f(\cdot)$ , and  $\tilde{\epsilon}_i$  is the residual based on the least squares  
 258 local linear estimate. By the adaptive nonparametric regression result of Linton and Xiao  
 259 (2007), the asymptotic result of  $\hat{\boldsymbol{\theta}}$  in (4.1) is the same whether the true density  $f(\cdot)$  is  
 260 used or not. Therefore, we will mainly show the existence and asymptotic distribution  
 261 of  $\hat{\boldsymbol{\theta}}$  assuming  $f(\cdot)$  is known.

We will first prove that with probability approaching 1, there exists a consistent local  
 maximizer  $\hat{\boldsymbol{\theta}} = (\hat{b}_1, \dots, \hat{b}_d, \hat{c}_1, \dots, \hat{c}_d)^T$  of (2.6) such that

$$\mathbf{H}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}) = O_p\{(nh)^{-1/2} + h^2\}.$$

262 Then we establish the asymptotic distributions for such consistent estimate.

Denote  $\boldsymbol{\theta}^* = \mathbf{H}\boldsymbol{\theta}$ ,  $\mathbf{x}_i^* = (x_{i1}, x_{i2}, \dots, x_{id}, (\frac{u_i - u_0}{h})x_{i1}, \dots, (\frac{u_i - u_0}{h})x_{id})^T$ ,  $K_i = K_h(u_i - u_0)$ ,  
 $R(u_i, \mathbf{x}_i) = \sum_{j=1}^d g_j(u_i)x_{ij} - \sum_{j=1}^d [b_j + c_j(u_i - u_0)]x_{ij}$ , and  $a_n = (nh)^{-1/2} + h^2$ . Let  
 $\rho(\cdot) = \log f(\cdot)$ , we have the objective function

$$L(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n K_i \rho(y_i - \boldsymbol{\theta}^{*T} \mathbf{x}_i^*) = L(\boldsymbol{\theta}^*).$$

It is sufficient to show that for any given  $\eta > 0$ , there exists a large constant  $c$  such that

$$P \left\{ \sup_{\|\mu\|=c} L(\boldsymbol{\theta}^* + a_n \mu) < L(\boldsymbol{\theta}^*) \right\} \geq 1 - \eta,$$

where  $\mu$  has the same dimension as  $\boldsymbol{\theta}$ ,  $a_n$  is the convergence rate. By using Taylor expansion, it follows that

$$\begin{aligned} L(\boldsymbol{\theta}^* + a_n \mu) - L(\boldsymbol{\theta}^*) &= \frac{1}{n} \sum_{i=1}^n K_i \{ \rho(\epsilon_i + R(u_i, \mathbf{x}_i) - a_n \mu^T \mathbf{x}_i^*) - \rho(\epsilon_i + R(u_i, \mathbf{x}_i)) \} \\ &= -\frac{1}{n} \sum_{i=1}^n K_i \rho'(\epsilon_i + R(u_i, \mathbf{x}_i)) a_n \mu^T \mathbf{x}_i^* + \frac{1}{2n} \sum_{i=1}^n K_i \rho''(\epsilon_i + R(u_i, \mathbf{x}_i)) a_n^2 (\mu^T \mathbf{x}_i^*)^2 \\ &\quad - \frac{1}{6n} \sum_{i=1}^n K_i \rho'''(z_i) a_n^3 (\mu^T \mathbf{x}_i^*)^3 \\ &\triangleq I_1 + I_2 + I_3, \end{aligned}$$

where  $z_i$  is a value between  $\epsilon_i + R(u_i, \mathbf{x}_i) - a_n \mu^T \mathbf{x}_i^*$  and  $\epsilon_i + R(u_i, \mathbf{x}_i)$ . For  $I_1 = -\frac{1}{n} \sum_{i=1}^n K_i \rho'(\epsilon_i + R(u_i, \mathbf{x}_i)) a_n \mu^T \mathbf{x}_i^*$ ,  $E(I_1) = -E \left[ K_i \rho'(\epsilon_i + R(u_i, \mathbf{x}_i)) a_n \mu^T \mathbf{x}_i^* \right]$ . We have,

$$\rho'(\epsilon_i + R(u_i, \mathbf{x}_i)) \approx \rho'(\epsilon_i) + \rho''(\epsilon_i) R(u_i, \mathbf{x}_i) + \frac{1}{2} \rho'''(\epsilon_i) R^2(u_i, \mathbf{x}_i).$$

Based on the assumption that  $\epsilon$  is independent of  $u$  and  $\mathbf{x}$ , and  $E[\rho'(\epsilon_i)] = 0$ , we have

$$E(I_1) \approx -a_n E \left\{ K_i \left[ \rho''(\epsilon_i) R(u_i, \mathbf{x}_i) + \frac{1}{2} \rho'''(\epsilon_i) R^2(u_i, \mathbf{x}_i) \right] \mu^T \mathbf{x}_i^* \right\}.$$

Since

$$\begin{aligned}
R(u_i, \mathbf{x}_i) &= \sum_{j=1}^d g_j(u_i) x_{ij} - \sum_{j=1}^d [b_j + c_j(u_i - u_0)] x_{ij} \\
&= \sum_{j=1}^d \left[ \sum_{m=2}^{\infty} \frac{1}{m!} g_j^{(m)}(u_0) (u_i - u_0)^m \right] x_{ij} \\
&= O_p(h^2),
\end{aligned}$$

then  $\frac{1}{2} \rho'''(\epsilon_i) R^2(u_i, \mathbf{x}_i) = [O_p(h^2)]^2 = O_p(h^4)$ , which is a smaller order than  $\rho''(\epsilon_i) R(u_i, \mathbf{x}_i)$ .

Thus,

$$E(I_1) \approx -a_n E \left\{ K_i \rho''(\epsilon_i) R(u_i, \mathbf{x}_i) \mu^T \mathbf{x}_i^* \right\} = -a_n E \left[ \rho''(\epsilon_i) \right] E \left[ K_i R(u_i, \mathbf{x}_i) \mu^T \mathbf{x}_i^* \right].$$

Since  $\delta_1 = E \left\{ \rho''(\epsilon_i) \right\}$ , then

$$E(I_1) \approx -a_n \delta_1 E \left[ K_i R(u_i, \mathbf{x}_i) \mu^T \mathbf{x}_i^* \right] = -a_n \delta_1 E \left\{ E \left\{ R(u_i, \mathbf{x}_i) \mu^T \mathbf{x}_i^* | u_i \right\} K_i \right\}.$$

By  $\mu^T \mathbf{x}_i^* \leq \|\mu\| \cdot \|\mathbf{x}_i^*\| = c \|\mathbf{x}_i^*\|$ , we have  $E(I_1) = O(a_n c h^2)$ .

$$\text{var}(I_1) = \frac{1}{n} \text{var} \left\{ K_i \rho'(\epsilon_i + R(u_i, \mathbf{x}_i)) a_n \mu^T \mathbf{x}_i^* \right\} = \frac{1}{n} \{ E(A^2) - [E(A)]^2 \},$$

where  $A = K_i \rho'(\epsilon_i + R(u_i, \mathbf{x}_i)) a_n \mu^T \mathbf{x}_i^*$ . Since  $\delta_2 = E \left\{ \rho'(\epsilon_i)^2 \right\}$ , then

$$\begin{aligned}
E(A^2) &= E \left\{ K_i^2 \rho'(\epsilon_i + R(u_i, \mathbf{x}_i))^2 a_n^2 (\mu^T \mathbf{x}_i^*)^2 \right\} \\
&\approx a_n^2 E \left\{ K_i^2 \rho'(\epsilon_i)^2 (\mu^T \mathbf{x}_i^*)^2 \right\} \\
&= a_n^2 \delta_2 E \left\{ E \left\{ (\mu^T \mathbf{x}_i^*)^2 | u_i \right\} K_i^2 \right\} \\
&= O \left( a_n^2 c^2 \frac{1}{h} \right).
\end{aligned}$$

Note that  $[E(A)]^2 = [O(a_n ch^2)]^2 \ll E(A^2)$ , then  $\text{var}(I_1) \approx \frac{1}{n}E(A^2) = O(a_n^2 c^2 \frac{1}{nh})$ . Hence,  $I_1 = E(I_1) + O_p(\sqrt{\text{var}(I_1)}) = O_p(a_n ch^2) + O_p\left(\sqrt{a_n^2 c^2 \frac{1}{nh}}\right) = O_p(ca_n^2)$ . For  $I_2 = \frac{1}{2n} \sum_{i=1}^n K_i \rho''(\epsilon_i + R(u_i, \mathbf{x}_i)) a_n^2 (\mu^T \mathbf{x}_i^*)^2$ ,

$$\begin{aligned} E(I_2) &= \frac{1}{2} a_n^2 E \left\{ K_i \rho''(\epsilon_i + R(u_i, \mathbf{x}_i)) (\mu^T \mathbf{x}_i^*)^2 \right\} \\ &= \frac{1}{2} a_n^2 E \left\{ \rho''(\epsilon_i) K_i (\mu^T \mathbf{x}_i^*)^2 \right\} (1 + o(1)) \\ &= \frac{1}{2} a_n^2 \delta_1 E \left\{ E \left\{ \mu^T \mathbf{x}_i^* \mathbf{x}_i^{*T} \mu | u_i \right\} K_i \right\} (1 + o(1)) \\ &= \frac{1}{2} a_n^2 \delta_1 \mu^T E \left\{ E \left\{ \mathbf{x}_i^* \mathbf{x}_i^{*T} | u_i \right\} K_i \right\} \mu (1 + o(1)). \end{aligned}$$

Note that  $\mathbf{x}_i^* \mathbf{x}_i^{*T} = \left( x_{ij} x_{ik} \left( \frac{u_i - u_0}{h} \right)^l \right)_{1 \leq j, k \leq d, l=0,1,2}$  and  $\Gamma_{jk}(u_i) = E(x_{ij} x_{ik} | u_i)$  for  $1 \leq j, k \leq d$ , then

$$\begin{aligned} E \left\{ E \left\{ x_{ij} x_{ik} \left( \frac{u_i - u_0}{h} \right)^l | u_i \right\} K_i \right\} &= E \left\{ E(x_{ij} x_{ik} | u_i) \left( \frac{u_i - u_0}{h} \right)^l K_i \right\} \\ &= E \left\{ \Gamma_{jk}(u_i) \left( \frac{u_i - u_0}{h} \right)^l K_i \right\}. \end{aligned}$$

By using Taylor expansion, we obtain

$$\begin{aligned} E \left\{ E \left\{ x_{ij} x_{ik} \left( \frac{u_i - u_0}{h} \right)^l | u_i \right\} K_i \right\} &= \frac{1}{h} \int \Gamma_{jk}(u_i) \left( \frac{u_i - u_0}{h} \right)^l K \left( \frac{u_i - u_0}{h} \right) q(u_i) du_i \\ &= q(u_0) \Gamma_{jk}(u_0) \int t^l K(t) dt (1 + o(1)). \end{aligned}$$

So we have

$$E(I_2) = \frac{1}{2} a_n^2 \delta_1 q(u_0) \mu^T \mathbf{S} \mu (1 + o(1)),$$



where  $\mathbf{S} = \begin{pmatrix} 1 & 0 \\ 0 & \mu_2 \end{pmatrix} \otimes \Gamma(u_0)$  is a  $2d \times 2d$  matrix. Thus,

$$\mathbb{E}(I_2) = O(a_n^2 \delta_1 q(u_0) \mu^T \mathbf{S} \mu)$$

and

$$\begin{aligned} \text{var}(I_2) &= \frac{a_n^4}{4n} \text{var} \left[ \rho''(\epsilon_i + R(u_i, \mathbf{x}_i)) K_i(\mu^T \mathbf{x}_i^*)^2 \right] \\ &= \frac{a_n^4}{4n} \{ \mathbb{E}(B^2) - [\mathbb{E}(B)]^2 \}, \end{aligned}$$

where  $B = \rho''(\epsilon_i + R(u_i, \mathbf{x}_i)) K_i(\mu^T \mathbf{x}_i^*)^2$ . Let  $\delta_3 = \mathbb{E}(\rho''(\epsilon_i)^2)$ , then

$$\begin{aligned} \mathbb{E}(B^2) &= \mathbb{E} \left\{ \rho''(\epsilon_i + R(u_i, \mathbf{x}_i))^2 K_i^2(\mu^T \mathbf{x}_i^*)^4 \right\} \\ &\approx \mathbb{E} \left\{ \rho''(\epsilon_i)^2 K_i^2(\mu^T \mathbf{x}_i^*)^4 \right\} \\ &= \delta_3 \mathbb{E} \left\{ K_i^2(\mu^T \mathbf{x}_i^*)^4 \right\} \\ &= O\left(\frac{1}{h}\right). \end{aligned}$$

Note that  $[\mathbb{E}(B)]^2 = [O(1)]^2 = O(1) \ll \mathbb{E}(B^2)$ , so  $\text{var}(I_2) = O\left(\frac{a_n^4}{nh}\right)$ . Based on the result  $I_2 = \mathbb{E}(I_2) + O_p(\sqrt{\text{var}(I_2)})$  and the assumption  $nh \rightarrow \infty$ , it follows that

$$I_2 = a_n^2 \delta_1 q(u_0) \mu^T \mathbf{S} \mu (1 + o_p(1)).$$

263 Similarly,  $I_3 = -\frac{1}{6n} \sum_{i=1}^n K_i \rho'''(z_i) a_n^3 (\mu^T \mathbf{x}_i^*)^3 = O_p(a_n^3)$ .

264 Assume  $\delta_1 < 0$ . Noticing that  $\mathbf{S}$  is a positive matrix,  $\|\mu\| = c$ , we can choose  $c$   
265 large enough such that  $I_2$  dominates both  $I_1$  and  $I_3$  with probability at least  $1 - \eta$ .

266 Thus  $P \left\{ \sup_{\|\mu\|=c} L(\boldsymbol{\theta}^* + a_n \mu) < L(\boldsymbol{\theta}^*) \right\} \geq 1 - \eta$ . Hence with probability approaching 1,

267 there exists a local maximizer  $\hat{\boldsymbol{\theta}}^*$  such that  $\|\hat{\boldsymbol{\theta}}^* - \boldsymbol{\theta}^*\| \leq a_n c$ , where  $a_n = (nh)^{-1/2} + h^2$ .  
 268 Based on the definition of  $\boldsymbol{\theta}^*$ , we can get, with probability approaching 1,  $H(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}) =$   
 269  $O_p((nh)^{-1/2} + h^2)$ .  $\square$

## 270 Proof of Theorem 2.2

Now we provide the asymptotic distribution for such consistent estimate. Since  $\hat{\boldsymbol{\theta}}$  maximizes  $L(\boldsymbol{\theta})$ , then  $L'(\hat{\boldsymbol{\theta}}) = 0$ . By Taylor expansion,

$$0 = L'(\hat{\boldsymbol{\theta}}) = L'(\boldsymbol{\theta}_0) + L''(\boldsymbol{\theta}_0)(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0) + \frac{1}{2}L'''(\tilde{\boldsymbol{\theta}})(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0)^2,$$

where  $\tilde{\boldsymbol{\theta}}$  is a value between  $\hat{\boldsymbol{\theta}}$  and  $\boldsymbol{\theta}_0$ . Then  $\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0 = -[L''(\boldsymbol{\theta}_0)]^{-1}L'(\boldsymbol{\theta}_0)(1 + o_p(1))$ . Since  $L(\boldsymbol{\theta}) = L(\boldsymbol{\theta}^*) = \frac{1}{n} \sum_{i=1}^n K_i \rho(y_i - \boldsymbol{\theta}^{*T} \mathbf{x}_i^*)$  and  $y_i - \boldsymbol{\theta}^{*T} \mathbf{x}_i^* = \epsilon_i + R(u_i, \mathbf{x}_i)$ , then  $L''(\boldsymbol{\theta}^*) = \frac{1}{n} \sum_{i=1}^n K_i \rho''(\epsilon_i + R(u_i, \mathbf{x}_i)) \mathbf{x}_i^* \mathbf{x}_i^{*T}$ . We have the following expectation,

$$\begin{aligned} \mathbb{E}[L''(\boldsymbol{\theta}^*)] &= \mathbb{E} \left\{ \rho''(\epsilon_i + R(u_i, \mathbf{x}_i)) K_i \mathbf{x}_i^* \mathbf{x}_i^{*T} \right\} \\ &\approx \mathbb{E} \left\{ \rho''(\epsilon_i) K_i \mathbf{x}_i^* \mathbf{x}_i^{*T} \right\} \\ &= \delta_1 \mathbb{E} \left\{ \mathbb{E} \left\{ \mathbf{x}_i^* \mathbf{x}_i^{*T} \mid u_i \right\} K_i \right\} \\ &= \delta_1 q(u_0) \mathbf{S} (1 + o(1)). \end{aligned}$$

Throughout this article, we consider the element-wise variance of a matrix,

$$\text{var}[L''(\boldsymbol{\theta}^*)] = \frac{1}{n} \text{var} \left\{ K_i \rho''(\epsilon_i + R(u_i, \mathbf{x}_i)) \mathbf{x}_i^* \mathbf{x}_i^{*T} \right\} = O \left( \frac{1}{nh} \right).$$

Based on the result  $L''(\boldsymbol{\theta}^*) = \text{E}[L''(\boldsymbol{\theta}^*)] + O_p(\sqrt{\text{var}[L''(\boldsymbol{\theta}^*)]})$  and the assumption  $nh \rightarrow \infty$ , it follows that

$$L''(\boldsymbol{\theta}^*) = \delta_1 q(u_0) \mathbf{S}(1 + o_p(1)).$$

For  $L'(\boldsymbol{\theta}^*)$ , we can divide it into two parts.

$$\begin{aligned} L'(\boldsymbol{\theta}^*) &= -\frac{1}{n} \sum_{i=1}^n K_i \rho'(\epsilon_i + R(u_i, \mathbf{x}_i)) \mathbf{x}_i^* \\ &\approx -\frac{1}{n} \sum_{i=1}^n K_i \rho'(\epsilon_i) \mathbf{x}_i^* - \frac{1}{n} \sum_{i=1}^n K_i \rho''(\epsilon_i) R(u_i, \mathbf{x}_i) \mathbf{x}_i^* \\ &\triangleq -\mathbf{w}_n - \boldsymbol{\nu}_n. \end{aligned}$$

The asymptotic result is determined by  $\mathbf{w}_n$ . In order to find the order of  $\boldsymbol{\nu}_n$ , we compute the following things.

$$\text{E}(\boldsymbol{\nu}_n) = \text{E} \left[ K_i \rho''(\epsilon_i) R(u_i, \mathbf{x}_i) \mathbf{x}_i^* \right] = \delta_1 \text{E} \left\{ \text{E} \left\{ R(u_i, \mathbf{x}_i) \mathbf{x}_i^* \mid u_i \right\} K_i \right\}.$$

Since  $g_j'''(\cdot)$  is bounded, then we have

$$R(u_i, \mathbf{x}_i) = \sum_{j=1}^d \left\{ \sum_{m=2}^{\infty} \frac{1}{m!} g_j^{(m)}(u_0) (u_i - u_0)^m \right\} x_{ij} = \sum_{j=1}^d \frac{1}{2} g_j''(u_0) (u_i - u_0)^2 x_{ij} (1 + o_p(1)).$$

By  $\mathbf{x}_i^* = (x_{i1}, \dots, x_{id}, (\frac{u_i - u_0}{h})x_{i1}, \dots, (\frac{u_i - u_0}{h})x_{id})^T$ ,

$$R(u_i, \mathbf{x}_i) \mathbf{x}_i^* \approx \left[ \left( \frac{(u_i - u_0)^2}{2} \left\{ \sum_{j=1}^d g_j''(u_0) x_{ij} \right\} x_{ik} \right)_{1 \leq k \leq d}, \left( \frac{(u_i - u_0)^3}{2h} \left\{ \sum_{j=1}^d g_j''(u_0) x_{ij} \right\} x_{ik} \right)_{1 \leq k \leq d} \right]_{2d \times 1}^T.$$

Since

$$\begin{aligned}
& \mathbb{E} \left\{ \mathbb{E} \left\{ \left[ \sum_{j=1}^d g_j''(u_0) x_{ij} \right] x_{ik} | u_i \right\} \frac{(u_i - u_0)^2}{2} K_i \right\} \\
&= \mathbb{E} \left\{ \sum_{j=1}^d g_j''(u_0) \mathbb{E}(x_{ij} x_{ik} | u_i) \frac{(u_i - u_0)^2}{2} K_i \right\} \\
&= \mathbb{E} \left\{ \sum_{j=1}^d g_j''(u_0) \Gamma_{jk}(u_i) \frac{(u_i - u_0)^2}{2} K_i \right\} \\
&= \sum_{j=1}^d g_j''(u_0) \mathbb{E} \left\{ \Gamma_{jk}(u_i) \frac{(u_i - u_0)^2}{2} K_i \right\} \\
&= \sum_{j=1}^d g_j''(u_0) \frac{1}{h} \int \Gamma_{jk}(u_i) \frac{(u_i - u_0)^2}{2} K\left(\frac{u_i - u_0}{h}\right) q(u_i) du_i \\
&= \frac{h^2}{2} q(u_0) \sum_{j=1}^d g_j''(u_0) \Gamma_{jk}(u_0) \int t^2 K(t) dt (1 + o(1))
\end{aligned}$$

and

$$\begin{aligned}
& \mathbb{E} \left\{ \mathbb{E} \left\{ \left[ \sum_{j=1}^d g_j''(u_0) x_{ij} \right] x_{ik} | u_i \right\} \frac{(u_i - u_0)^3}{2h} K_i \right\} \\
&= \mathbb{E} \left\{ \sum_{j=1}^d g_j''(u_0) \Gamma_{jk}(u_i) \frac{(u_i - u_0)^3}{2h} K_i \right\} \\
&= \sum_{j=1}^d g_j''(u_0) \frac{1}{2h} \mathbb{E} \left\{ \Gamma_{jk}(u_i) (u_i - u_0)^3 K_i \right\} \\
&= \frac{h^2}{2} q(u_0) \sum_{j=1}^d g_j''(u_0) \Gamma_{jk}(u_0) \int t^3 K(t) dt (1 + o(1)),
\end{aligned}$$

then

$$\mathbb{E}(\boldsymbol{\nu}_n) = \delta_1 q(u_0) \frac{h^2}{2} \sum_{j=1}^d g_j''(u_0) \boldsymbol{\psi}_j (1 + o(1)),$$

where  $\boldsymbol{\psi}_j = \begin{pmatrix} \mu_2 \\ \mu_3 \end{pmatrix} \otimes (\Gamma_{jk}(u_0))_{1 \leq k \leq d}^T$  is a  $2d \times 1$  vector for  $j = 1, \dots, d$ . Since  $\text{var}(\boldsymbol{\nu}_n) = \text{var} \{ K_i \rho_j''(\epsilon_i) R(u_i, \mathbf{x}_i) \mathbf{x}_i^* \} / n = O(h^3/n)$ , then based on the result  $\boldsymbol{\nu}_n = \mathbb{E}(\boldsymbol{\nu}_n) + O_p(\sqrt{\text{var}(\boldsymbol{\nu}_n)})$

and the assumption  $nh \rightarrow \infty$ , it follows that

$$\boldsymbol{\nu}_n = \delta_1 q(u_0) \frac{h^2}{2} \sum_{j=1}^d g_j''(u_0) \boldsymbol{\psi}_j (1 + o_p(1)).$$

Then

$$\begin{aligned} \hat{\boldsymbol{\theta}}^* - \boldsymbol{\theta}^* &= -[L''(\boldsymbol{\theta}^*)]^{-1} L'(\boldsymbol{\theta}^*) (1 + o_p(1)) \\ &= -[\delta_1 q(u_0) \mathbf{S}]^{-1} (-\mathbf{w}_n - \boldsymbol{\nu}_n) (1 + o_p(1)) \\ &= \frac{\mathbf{S}^{-1} \mathbf{w}_n}{\delta_1 q(u_0)} (1 + o_p(1)) + \mathbf{S}^{-1} \frac{h^2}{2} \sum_{j=1}^d g_j''(u_0) \boldsymbol{\psi}_j (1 + o_p(1)). \end{aligned} \quad (4.2)$$

Based on the assumption  $E[\rho'(\epsilon_i)] = 0$ , we can easily get  $E(\mathbf{w}_n) = 0$ .

$$\text{var}(\mathbf{w}_n) = \frac{1}{n} \text{var} \left\{ K_i \rho'(\epsilon_i) \mathbf{x}_i^* \right\} = \frac{1}{n} E \left\{ K_i^2 \rho'(\epsilon_i)^2 \mathbf{x}_i^* \mathbf{x}_i^{*T} \right\} = \frac{1}{n} \delta_2 E \left\{ E \left\{ \mathbf{x}_i^* \mathbf{x}_i^{*T} | u_i \right\} K_i^2 \right\}.$$

Since  $\mathbf{x}_i^* \mathbf{x}_i^{*T} = \left( x_{ij} x_{ik} \left( \frac{u_i - u_0}{h} \right)^l \right)_{1 \leq j, k \leq d, l=0,1,2}$  and

$$\begin{aligned} E \left\{ E \left\{ x_{ij} x_{ik} \left( \frac{u_i - u_0}{h} \right)^l | u_i \right\} K_i^2 \right\} &= E \left\{ E \left\{ x_{ij} x_{ik} | u_i \right\} \left( \frac{u_i - u_0}{h} \right)^l K_i^2 \right\} \\ &= E \left\{ \Gamma_{jk}(u_i) \left( \frac{u_i - u_0}{h} \right)^l K_i^2 \right\} \\ &= \frac{1}{h} q(u_0) \Gamma_{jk}(u_0) \int t^l K^2(t) dt (1 + o(1)), \end{aligned}$$

then

$$E \left\{ E \left\{ \mathbf{x}_i^* \mathbf{x}_i^{*T} | u_i \right\} K_i^2 \right\} = \frac{1}{h} q(u_0) \boldsymbol{\Lambda} (1 + o(1)),$$

where  $\boldsymbol{\Lambda} = \begin{pmatrix} \nu_0 & \nu_1 \\ \nu_1 & \nu_2 \end{pmatrix} \otimes \boldsymbol{\Gamma}(u_0)$  is a  $2d \times 2d$  matrix. So  $\text{var}(\mathbf{w}_n) = \frac{1}{nh} \delta_2 q(u_0) \boldsymbol{\Lambda} (1 + o(1))$ .

We next use the Lyapunov central limit theorem to obtain the asymptotic distribution

of  $\mathbf{w}_n$ . The Lyapunov conditions are checked as follows. For any unit vector  $\mathbf{d} \in \mathbb{R}^{2d}$ , let  $\mathbf{d}^T \mathbf{w}_n = \sum_{i=1}^n \xi_i$ , where  $\xi_i = \frac{1}{n} K_i \rho'(\epsilon_i) \mathbf{d}^T \mathbf{x}_i^*$ . Since

$$E(\xi_i^2) = E \left\{ \frac{1}{n^2} K_i^2 \rho'(\epsilon_i)^2 \mathbf{d}^T \mathbf{x}_i^* \mathbf{x}_i^{*T} \mathbf{d} \right\} = \frac{1}{n^2} \delta_2 \mathbf{d}^T E \{ K_i^2 \mathbf{x}_i^* \mathbf{x}_i^{*T} \} \mathbf{d} = \frac{1}{n^2 h} \delta_2 q(u_0) \mathbf{d}^T \mathbf{\Lambda} \mathbf{d} (1 + o(1)),$$

then  $o \left( \left( \sum_{i=1}^n E |\xi_i|^2 \right)^3 \right) = o \left( \left( \frac{1}{nh} \right)^3 \right)$ . Let  $\delta_4 = E \{ \rho'(\epsilon_i)^3 \}$ , then

$$E(\xi_i^3) = E \left\{ \frac{1}{n^3} K_i^3 \rho'(\epsilon_i)^3 (\mathbf{d}^T \mathbf{x}_i^*)^3 \right\} = \frac{1}{n^3} \delta_3 E \{ K_i^3 (\mathbf{d}^T \mathbf{x}_i^*)^3 \} = O \left( \frac{1}{n^3 h^2} \right).$$

So  $\left( \sum_{i=1}^n E |\xi_i|^3 \right)^2 = O \left( \left( \frac{1}{n^2 h^2} \right)^2 \right)$ . Since  $\left( \frac{1}{n^2 h^2} \right)^2 (nh)^3 = \frac{1}{nh} \rightarrow 0$ , then  $\left( \frac{1}{n^2 h^2} \right)^2 = o \left( \left( \frac{1}{nh} \right)^3 \right)$ , which is equivalent to  $\left( \sum_{i=1}^n E |\xi_i|^3 \right)^2 = o \left( \left( \sum_{i=1}^n E |\xi_i|^2 \right)^3 \right)$ . Based on Lyapunov Central Limit Theorem,

$$\frac{\mathbf{w}_n}{\sqrt{\text{var}(\mathbf{w}_n)}} \xrightarrow{D} N(\mathbf{0}_{2d}, \mathbf{I}_{2d}),$$

where  $\mathbf{0}_{2d}$  is a  $2d \times 1$  vector with each entry being 0;  $\mathbf{I}_{2d}$  is a  $2d \times 2d$  identity matrix. Previously, we already computed that  $\text{var}(\mathbf{w}_n) = \frac{1}{nh} \delta_2 q(u_0) \mathbf{\Lambda} (1 + o(1))$ , by Slutsky's Theorem,

$$\sqrt{nh} \mathbf{w}_n \xrightarrow{D} N(\mathbf{0}_{2d}, \delta_2 q(u_0) \mathbf{\Lambda}).$$

Based on (4.2), we have the following result

$$\sqrt{nh} \left\{ \mathbf{H}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}) - \mathbf{S}^{-1} \frac{h^2}{2} \sum_{j=1}^d g_j''(u_0) \boldsymbol{\psi}_j (1 + o_p(1)) \right\} \xrightarrow{D} N(\mathbf{0}_{2d}, \delta_1^{-2} \delta_2 q(u_0)^{-1} \mathbf{S}^{-1} \mathbf{\Lambda} \mathbf{S}^{-1}).$$

## 271 References

- 272 Beran, R. (1974). Asymptotic efficient adaptive rank estimates in location models. *The*  
 273 *Annals of Statistics*, 2, 63-74.

- 274 Bickel, P. J. (1982). On adaptive estimation. *The Annals of Statistics*, 10, 647-671.
- 275 Cai, Z., Fan, J. and Li, R. (2000). Efficient estimation and inferences for varying-  
276 coefficient models. *Journal of the American Statistical Association*, 95, 888-902.
- 277 Chiang, C-T., Rice, J. A. and Wu, C. O. (2001). Smoothing spline estimation for vary-  
278 ing coefficient models with repeatedly measured dependent variable. *Journal of the*  
279 *American Statistical Association*, 96, 605-619.
- 280 Cleveland, W. S., Grosse, E. and Shyu, W. M. (1991). Local regression models. In *Statis-*  
281 *tical Models in S* (J.M. Chambers and T.J. Hastie, eds.), 309-376. Wadsworth/Brooks-  
282 Cole, Pacific Grove, CA.
- 283 Drost, F. C. and Klaassen, C. A. J. (1997). Efficient estimation in semiparametric  
284 GRACH models. *Journal of Econometrics*, 81, 193-221.
- 285 Eubank, R. L., Huang, C. F., Maldonado, Y. M., Wang, N., Wang, S. and Buchanan, R.  
286 J. (2004). Smoothing spline estimation in varying-coefficient models. *Journal of the*  
287 *Royal Statistical Society*, Ser. B, 66, 653-667.
- 288 Fan, J. (1993). Local linear regression smoothers and their minimax efficiencies. *The*  
289 *Annals of Statistics*, 21, 196-216.
- 290 Fan, J., Farmen, M. and Gijbels, I. (1998). Local maximum likelihood estimation and  
291 inference. *Journal of the Royal Statistical Society*, B, 60, 591-608.
- 292 Fan J. and Gijbels, I. (1992). Variable Bandwidth and Local Linear Regression  
293 Smoothers. *The Annals of Statistics*, 20, 2008-2036.
- 294 Fan, J. and Huang, T. (2005). Profile likelihood inferences on semiparametric varying-  
295 coefficient partially linear models. *Bernoulli*, 11, 1031-1057.

- 296 Fan, J. and Zhang, W. (1999). Statistical estimation in varying coefficient models. *The*  
297 *Annals of Statistics*, 27, 1491-1518.
- 298 Fan, J. and Zhang, J. T. (2000). Two-step estimation of functional linear models with  
299 applications to longitudinal data. *Journal of the Royal Statistical Society*, Ser. B, 62,  
300 303-322.
- 301 Fan, J. and Zhang, W. (2008). Statistical methods with varying coefficient models.  
302 *Statistics and Its Interface*, 1(1), 179-195.
- 303 Fan, J., Zhang, C. and Zhang, J. (2001). Generalized likelihood ratio statistics and Wilks  
304 phenomenon. *The Annals of Statistics*, 29, 153-193.
- 305 Gilley, O. W. and Pace, R. K. (1996). On the Harrison and Rubinfeld data. *Journal of*  
306 *Environmental Economics and Management*, 31, 403-405.
- 307 Hastie, T. J. and Tibshirani, R. J. (1993). Varying-Coefficient Models. *Journal of the*  
308 *Royal Statistical Society*, Ser. B, 55, 757-796.
- 309 Hodgson, D. J. (1998). Adaptive estimation of cointegrating regressions with ARMA  
310 errors. *Journal of Econometrics*, 85, 231-267.
- 311 Hoover, D. R., Rice, J. A., Wu, C. O. and Yang, L. P. (1998). Nonparametric smoothing  
312 estimates of time-varying coefficient models with longitudinal data. *Biometrika*, 85,  
313 809-822.
- 314 Huang, J. Z., Wu, C. O. and Zhou, L. (2002). Varying-coefficient models and basis  
315 function approximations for the analysis of repeated measurements. *Biometrika*, 89,  
316 111-128.
- 317 Huang, J. Z. and Shen, H. (2004). Functional coefficient regression models for non-linear



- 318 time series: A polynomial spline approach. *Scandinavian Journal of Statistics*, 31,  
319 515-534.
- 320 Huang, J. Z., Wu, C. O. and Zhou, L. (2004). Polynomial spline estimation and inference  
321 for varying coefficient models with longitudinal data. *Statistica Sinica*, 14, 763-788.
- 322 Kauermann, G. and Tutz, G. (1999). On model diagnostics using varying coefficient  
323 models. *Biometrika*, 86, 119-128.
- 324 Linton, O. B. and Xiao, Z. (2007). A nonparametric regression estimator that adapts to  
325 error distribution of unknown form. *Econometric Theory*, 23, 371-413.
- 326 Manski, C. F. (1984). Adaptive estimation of non-linear regression models. *Econometric*  
327 *Reviews*, 3, 145-194.
- 328 Schick, A. (1993). On efficient estimation in regression models. *The Annals of Statistics*,  
329 21, 1486-1521.
- 330 Staniswalis, J. G. (1989). The kernel estimate of a regression function in likelihood based  
331 models. *Journal of the American Statistical Association*, 84, 276-283.
- 332 Steigerwald, D. G. (1992). Adaptive estimation in time series regression models. *Journal*  
333 *of Econometrics*, 54, 251-275.
- 334 Stone, C. J. (1975). Adaptive maximum likelihood estimators of a location parameter.  
335 *The Annals of Statistics*, 3, 267-284.
- 336 Sun, Y., Yan, H., Zhang, W. and Lu, Z. (2014). A semiparametric spatial dynamic  
337 model. *The Annals of Statistics*, 42, 700-727.
- 338 Wang, L., Kai, B. and Li, R. (2009). Local rank inference for varying coefficient models.  
339 *Journal of the American Statistical Association*, 104, 1631-1645.

- 340 Wang, H. and Xia, Y. (2009). Shrinkage estimation of the varying coefficient model.  
341 *Journal of the American Statistical Association*, 104, 747-757.
- 342 Wang, Q. and Yao, W. (2012). An adaptive estimation of MAVE. *Journal of Multivariate*  
343 *Analysis*, 104, 88-100.
- 344 Wu, C. O., Chiang, C-T. and Hoover, D. R. (1998). Asymptotic Confidence Regions for  
345 Kernel Smoothing of a Varying Coefficient Model with Longitudinal Data. *Journal of*  
346 *the American Statistical Association*, 93, 1388-1402.
- 347 Xia, Y., Tong, H., Li, W. and Zhu, L. (2002). An adaptive estimation of dimension  
348 reduction space. *Journal of the Royal Statistical Society*, Ser. B, 64, 363-388.
- 349 Yao, W. (2013). A note on EM algorithm for mixture models. *Statistics and Probability*  
350 *Letters*, 83, 519-526.
- 351 Yao, W. and Zhao, Z. (2013). Kernel density based linear regression estimates. *Communi-*  
352 *cations in Statistics-Theory and Methods*, 42, 4499-4512.
- 353 Yuan, A. and De Gooijer, J. G. (2007). Semiparametric regression with kernel error  
354 model. *Scandinavian Journal of Statistics*, 34, 841-869.
- 355 Yuan, A. (2009). Semiparametric inference with kernel likelihood. *Journal of Nonpara-*  
356 *metric Statistics*, 21, 207-228.
- 357 Zhang, W. and Lee, S. Y. (2000). Variable bandwidth selection in varying-coefficient  
358 models. *Journal of Multivariate Analysis*, 74, 116-134.