

This is a repository copy of *A semiparametric spatial dynamic model*.

White Rose Research Online URL for this paper:

<https://eprints.whiterose.ac.uk/id/eprint/94971/>

Version: Published Version

---

**Article:**

Sun, Yan, Yan, Hongjia, Zhang, Wenyang orcid.org/0000-0001-8391-1122 et al. (1 more author) (2014) A semiparametric spatial dynamic model. *Annals of Statistics*. pp. 700-727. ISSN: 0090-5364

<https://doi.org/10.1214/13-AOS1201>

---

**Reuse**

Items deposited in White Rose Research Online are protected by copyright, with all rights reserved unless indicated otherwise. They may be downloaded and/or printed for private study, or other acts as permitted by national copyright laws. The publisher or other rights holders may allow further reproduction and re-use of the full text version. This is indicated by the licence information on the White Rose Research Online record for the item.

**Takedown**

If you consider content in White Rose Research Online to be in breach of UK law, please notify us by emailing [eprints@whiterose.ac.uk](mailto:eprints@whiterose.ac.uk) including the URL of the record and the reason for the withdrawal request.

## A SEMIPARAMETRIC SPATIAL DYNAMIC MODEL<sup>1</sup>

BY YAN SUN, HONGJIA YAN, WENYANG ZHANG AND ZUDI LU

*Shanghai University of Finance and Economics, The University of York,  
 The University of York and The University of Southampton*

Stimulated by the Boston house price data, in this paper, we propose a semiparametric spatial dynamic model, which extends the ordinary spatial autoregressive models to accommodate the effects of some covariates associated with the house price. A profile likelihood based estimation procedure is proposed. The asymptotic normality of the proposed estimators are derived. We also investigate how to identify the parametric/nonparametric components in the proposed semiparametric model. We show how many unknown parameters an unknown bivariate function amounts to, and propose an AIC/BIC of nonparametric version for model selection. Simulation studies are conducted to examine the performance of the proposed methods. The simulation results show our methods work very well. We finally apply the proposed methods to analyze the Boston house price data, which leads to some interesting findings.

**1. Introduction.** The Boston house price data is frequently used in literature to illustrate some new statistical methods. If we use  $y_i$  to denote the median value of owner-occupied homes at location  $s_i$ , a spatial autoregressive model for the data would be

$$(1.1) \quad y_i = \sum_{j \neq i} w_{ij} y_j + \varepsilon_i, \quad i = 1, \dots, n,$$

where  $w_{ij}$  is the impact of  $y_j$  on  $y_i$ . However, (1.1) is inadequate because it models  $y_i$  solely based on the median value prices,  $y_j$ , for  $j \neq i$ . It is better to incorporate the effects of some important covariates, such as the crime rate and accessibility to radial highways, into the model. Let  $X_i$ , a  $p$ -dimensional vector, be the vector of the covariates associated with  $y_i$ . A reasonable model to fit the data would be

$$(1.2) \quad y_i = \sum_{j \neq i} w_{ij} y_j + X_i^T \boldsymbol{\beta} + \varepsilon_i, \quad i = 1, \dots, n,$$

---

Received May 2013; revised September 2013.

<sup>1</sup>Supported by National Science Foundation of China (Grant 11271242), program for New Century Excellent Talents in China's University (NECT-10-0562), Key Laboratory of Mathematical Economics (SUF), Ministry of Education of China and the Singapore National Research Foundation under its Cooperative Basic Research Grant and administered by the Singapore Ministry of Health's National Medical Research Council (Grant No. NMRC/CBRG/0014/2012).

*MSC2010 subject classifications.* Primary 62G08; secondary 62G05, 62G20.

*Key words and phrases.* AIC/BIC, local linear modeling, profile likelihood, spatial interaction.

where  $w_{ij}$  and  $\beta$  are unknown. However, there are two problems with model (1.2): first, there are too many unknown parameters; second, the model has not taken into account the location effects of the impacts of the covariates—the impacts of some covariates may vary over location. To control the number of unknown parameters and take the location effects into account, we propose the following model to fit the data:

$$(1.3) \quad y_i = \alpha \sum_{j \neq i} w_{ij} y_j + X_i^T \beta(s_i) + \varepsilon_i, \quad i = 1, \dots, n,$$

where  $w_{ij}$  is a specified certain physical or economic distance,  $s_i$  is the location of the  $i$ th observation, which is a two-dimensional vector,  $\beta(\cdot) = (\beta_1(\cdot), \dots, \beta_p(\cdot))^T$ ,  $\varepsilon_i, i = 1, \dots, n$ , are i.i.d., and follow  $N(0, \sigma^2)$ ,  $\{X_i, i = 1, \dots, n\}$  is independent of  $\{\varepsilon_i, i = 1, \dots, n\}$ .  $\alpha, \sigma^2$  and  $\beta(\cdot)$  are unknown to be estimated. Model (1.3) is the model this paper is going to address. From now on,  $y_i$  is of course not necessarily the house price, it is a generic response variable. We will also see that the normality assumption imposed on  $\varepsilon_i$  is just for the description of the construction of the proposed estimation procedure. It is not necessary for the asymptotic properties of the proposed estimators.

In model (1.3), the spatial neighboring effect of  $y_j, j \neq i$ , on  $y_i$  is formulated through  $\alpha w_{ij}$ , where  $w_{ij}$  is a specified certain physical or economic distance, and  $\alpha$  is an unknown baseline of the spatial neighboring effect. Such method to define spatial neighboring effect is common; see Ord [12], Anselin [1], Su and Jin [13].

If there is no any condition imposed on the spatial neighboring effects, and the spatial neighboring effects are formulated as unknown  $w_{ij}, i = 1, \dots, n, j = 1, \dots, i-1, i+1, \dots, n$ , we would have  $(n-1)n$  unknown  $w_{ij}$ 's to estimate. In which case, it would be impossible to have consistent estimators of  $w_{ij}$ 's. However, if we impose some kind of sparsity on  $w_{ij}$ 's, by penalized maximum likelihood estimation, it is possible to construct consistent estimators of  $w_{ij}$ 's. However, that has gone beyond the scope of this paper although it is a promising research project.

Model (1.3) is a useful extension of spatial autoregressive models (Gao et al. [6]; Kelejian and Prucha [8]; Ord [12]; Su and Jin [13]) and varying coefficient models (Cheng et al. [2]; Fan and Zhang [4, 5]; Li and Zhang [10]; Sun et al. [15]; Zhang et al. [19, 20]; Wang and Xia [17]; and Tao and Xia [16]). One characteristic of model (1.3) is

$$E(\varepsilon_i | y_1, \dots, y_{i-1}, y_{i+1}, \dots, y_n) \neq 0$$

although  $E(\varepsilon_i) = 0$ , the standard least squares estimation will not work for (1.3). In this paper, based on the local linear modeling and profile likelihood idea, we will propose a local likelihood based estimation procedure for the unknown parameters and functions in (1.3) and derive the asymptotic properties of the obtained estimators.

In reality, some components of  $\beta(\cdot)$  in model (1.3) may be constant, and we do not know which components are functional, which are constant. Methodologically speaking, if mistakenly treating a constant component as functional, we would pay a price on the variance side of the obtained estimator; on the other hand, if mistakenly treating a functional component as constant, we would pay a price on the bias side of the obtained estimator. The identification of constant/functional components in  $\beta(\cdot)$  is imperative. From practical point of view, the identification of constant components is also of importance. For the data set we study in this paper,  $\beta(\cdot)$  can be interpreted as the vector of the impacts of the covariates concerned on the house price. The identification will reveal which covariates have location varying impacts on the house price, and which do not. This is apparently something of great interest. In this paper, we will show how many unknown parameters an unknown bivariate function amounts to, and propose an AIC/BIC of nonparametric version to identify the constant components of  $\beta(\cdot)$  in model (1.3).

The paper is organized as follows. We begin in Section 2 with a description of the estimation procedure for the proposed model (1.3). In Section 3, we show how many unknown parameters an unknown bivariate function amounts to, and propose an AIC/BIC of nonparametric version for model selection. Asymptotic properties of the proposed estimators are presented in Section 4. The performance of the proposed methods, including both estimation and model selection methods, is assessed by a simulation study in Section 5. In Section 6, we explore how the covariates, which are commonly found to be associated with house price, affect the median value of owner-occupied homes in Boston, and how the impacts of these covariates change over location based on the proposed model and estimation procedure.

Throughout this paper,  $\mathbf{0}_k$  is a  $k$ -dimensional vector with each component being 0,  $I_k$  is an identity matrix of size  $k$ ,  $U[0, 1]^2$  is a two-dimensional uniform distribution on  $[0, 1] \times [0, 1]$ .

**2. Estimation procedure.** Let  $w_{ii} = 0$ ,  $W = (w_{ij})$ ,  $Y = (y_1, \dots, y_n)^T$ ,  $A = I_n - \alpha W$  and  $\mathbf{m} = (X_1^T \beta(s_1), \dots, X_n^T \beta(s_n))^T$ . By simple calculations, we have that the conditional density function of  $Y$  given  $\mathbf{m}$  is  $N(A^{-1}\mathbf{m}, (A^T A)^{-1}\sigma^2)$ , which leads to the following log likelihood function

$$(2.1) \quad -\frac{n}{2} \log(2\pi) - n \log(\sigma) + \log(|A|) - \frac{1}{2\sigma^2} (AY - \mathbf{m})^T (AY - \mathbf{m}).$$

Our estimation is profile likelihood based. We first construct the estimator  $\tilde{\beta}(\cdot; \alpha)$  of  $\beta(\cdot)$  pretending  $\alpha$  is known, then let  $(\hat{\alpha}, \hat{\sigma}^2)$  maximize (2.1) with  $\beta(\cdot)$  being replaced by  $\tilde{\beta}(\cdot; \alpha)$ .  $\hat{\alpha}$  and  $\hat{\sigma}^2$  are our estimators of  $\alpha$  and  $\sigma^2$ , respectively. After the estimator of  $\alpha$  is obtained, the estimator of  $\beta(\cdot)$  is taken to be  $\tilde{\beta}(\cdot; \alpha)$  with  $\alpha$  and the bandwidth used being replaced by  $\hat{\alpha}$  and a slightly larger bandwidth, respectively. The details are as follows.

For any  $s = (u, v)^T$ , we denote  $(\partial \boldsymbol{\beta}(s)/\partial u, \partial \boldsymbol{\beta}(s)/\partial v)$  by  $\dot{\boldsymbol{\beta}}(s)$ , where  $\partial \boldsymbol{\beta}(s)/\partial u = (\partial \beta_1(s)/\partial u, \dots, \partial \beta_p(s)/\partial u)^T$ . We define  $\|s\| = (s^T s)^{1/2}$ .

For any given  $s$ , by the Taylor's expansion, we have

$$\boldsymbol{\beta}(s_i) \approx \boldsymbol{\beta}(s) + \dot{\boldsymbol{\beta}}(s)(s_i - s),$$

when  $s_i$  is in a small neighborhood of  $s$ , which leads to the following objective function for estimating  $\boldsymbol{\beta}(s)$ :

$$(2.2) \quad \sum_{i=1}^n (y_i^* - X_i^T \mathbf{a} - X_i^T \mathbf{B}(s_i - s))^2 K_h(\|s_i - s\|),$$

where  $y_i^*$  is the  $i$ th component of  $AY$ ,  $K_h(\cdot) = K(\cdot/h)/h^2$ ,  $K(\cdot)$  is a kernel function, and  $h$  is a bandwidth. Let  $(\hat{\mathbf{a}}, \hat{\mathbf{B}})$  minimise (2.2), the “estimator”  $\tilde{\boldsymbol{\beta}}(s; \alpha)$  of  $\boldsymbol{\beta}(s)$  is taken to be  $\hat{\mathbf{a}}$ . By simple calculations, we have

$$(2.3) \quad \tilde{\boldsymbol{\beta}}(s; \alpha) = \hat{\mathbf{a}} = (I_p, \mathbf{0}_{p \times 2p})(\mathcal{X}^T \mathcal{W} \mathcal{X})^{-1} \mathcal{X}^T \mathcal{W} AY,$$

where  $\mathbf{0}_{p \times q}$  is a matrix of size  $p \times q$  with each entry being 0, and

$$\mathcal{X} = \begin{pmatrix} X_1 & \cdots & X_n \\ X_1 \otimes (s_1 - s) & \cdots & X_n \otimes (s_n - s) \end{pmatrix}^T,$$

$$\mathcal{W} = \text{diag}(K_h(\|s_1 - s\|), \dots, K_h(\|s_n - s\|)).$$

Replacing  $\boldsymbol{\beta}(s_i)$  in (2.1) by  $\tilde{\boldsymbol{\beta}}(s_i; \alpha)$  and ignoring the constant term, we have the objective function for estimating  $\alpha$  and  $\sigma^2$

$$(2.4) \quad -n \log(\sigma) + \log(|A|) - \frac{1}{2\sigma^2} (AY - \tilde{\mathbf{m}})^T (AY - \tilde{\mathbf{m}}),$$

where  $\tilde{\mathbf{m}}$  is  $\mathbf{m}$  with  $\boldsymbol{\beta}(s_i)$  being replaced by  $\tilde{\boldsymbol{\beta}}(s_i; \alpha)$ . Let  $\alpha_i$ ,  $i = 1, \dots, n$ , be the eigenvalues of  $W$ ,

$$\tilde{\sigma}^2 = \frac{1}{n} (AY - \tilde{\mathbf{m}})^T (AY - \tilde{\mathbf{m}})$$

and  $(\hat{\alpha}, \hat{\sigma}^2)$  maximize (2.4). Noticing that  $|A| = \prod_{i=1}^n (1 - \alpha \alpha_i)$ , by simple calculations, we have  $\hat{\alpha}$  is the maximizer of

$$(2.5) \quad -n \log(\tilde{\sigma}) + \sum_{i=1}^n \log(|1 - \alpha \alpha_i|)$$

and  $\hat{\sigma}^2$  is  $\tilde{\sigma}^2$  with  $\alpha$  being replaced by  $\hat{\alpha}$ .

Note that the maximization of (2.5) is not difficult because it is a one-dimensional optimization problem, which can be solved using a grid point method.

The estimator  $\hat{\boldsymbol{\beta}}(\cdot) = (\hat{\beta}_1(\cdot), \dots, \hat{\beta}_p(\cdot))^T$  is  $\tilde{\boldsymbol{\beta}}(\cdot; \alpha)$  with  $\alpha$  being replaced by  $\hat{\alpha}$  and the bandwidth  $h$  by a slightly larger bandwidth  $h_1$ . The reason for replacing

the bandwidth  $h$  by a slightly larger number  $h_1$  is that the former bandwidth is appropriate for the estimation of constant parameters,  $\alpha$  and  $\alpha^2$ , and the latter is more appropriate for the estimation of functional parameters. Also, the estimators of constant parameters need a smaller bandwidth  $h$  in order to achieve the optimal rate of convergence.

In reality, some components of  $\beta(\cdot)$  may be constant. If a component of  $\beta(\cdot)$  is a constant, say  $\beta_1(\cdot) = \beta_1$ , we use the average of  $\hat{\beta}_1(s_i)$ ,  $i = 1, \dots, n$ , to estimate the constant  $\beta_1$ , that is,

$$\hat{\beta}_1 = \frac{1}{n} \sum_{i=1}^n \hat{\beta}_1(s_i).$$

How to identify the constant components of  $\beta(\cdot)$  will be addressed in the next section.

### 3. Identification of constant components.

**3.1. Criterion for identification.** As we mentioned before, some components of  $\beta(\cdot)$  in model (1.3) may be constant in reality, and to identify such constant components is of importance. In this paper, we appeal the AIC or BIC to identify the constant components. The AIC for (1.3), in which some components of  $\beta(\cdot)$  may be constant, is defined as follows:

$$(3.1) \quad \text{AIC} = n \log(\hat{\sigma}) - \log(|\hat{A}|) + \frac{1}{2\hat{\sigma}^2} (\hat{A}Y - \hat{\mathbf{m}})^T (\hat{A}Y - \hat{\mathbf{m}}) + \mathcal{K},$$

where  $\hat{A}$  and  $\hat{\mathbf{m}}$  are  $A$  and  $\mathbf{m}$  with the unknown parameters and functions being replaced by their estimators,  $\mathcal{K}$  is the number of unknown parameters in model (1.3). The BIC can be defined in a similar way.

Because there are unknown functions in model (1.3), the first hurdle in the calculation of AIC of model (1.3) is to find how many unknown constants an unknown bivariate function amounts to. In the following, based on the residual sum of squares of standard bivariate nonparametric regression model, we propose an ad hoc way to solve this problem.

Suppose we have the following standard bivariate nonparametric regression model:

$$(3.2) \quad \eta_i = g(s_i) + e_i, \quad i = 1, \dots, n,$$

where  $E(e_i) = 0$  and  $\text{var}(e_i) = \sigma_e^2$ . The residual sum of squares of (3.2) is

$$\text{RSS} = \sum_{i=1}^n \{\eta_i - \hat{g}(s_i)\}^2,$$

where  $\hat{g}(\cdot)$  is the local linear estimator of  $g(\cdot)$ . On the other hand,

$E(\text{RSS}/\sigma_e^2) = n -$  the number of unknown parameters in the regression function.

So, the number  $\mathcal{T}$  of unknown constants the unknown function  $g(\cdot)$  amounts to can be reasonably viewed as

$$\mathcal{T} = n - E(\text{RSS}/\sigma_e^2) = n - \sigma_e^{-2} E \left[ \sum_{i=1}^n \{\eta_i - \hat{g}(s_i)\}^2 \right].$$

To make  $\mathcal{T}$  more convenient to use, we derive the asymptotic form of  $\mathcal{T}$ . Let

$$\mathbf{S}_i = \begin{pmatrix} 1 & s_1^T - s_i^T \\ \vdots & \vdots \\ 1 & s_n^T - s_i^T \end{pmatrix}, \quad \boldsymbol{\eta} = \begin{pmatrix} \eta_1 \\ \vdots \\ \eta_n \end{pmatrix}, \quad \mathbf{e} = \begin{pmatrix} e_1 \\ \vdots \\ e_n \end{pmatrix}$$

and

$$\mathcal{W}_i = \text{diag}(K_h(u_1 - u_i)K_h(v_1 - v_i), \dots, K_h(u_n - u_i)K_h(v_n - v_i)),$$

we have

$$\hat{g}(s_i) = (1, 0, 0)(\mathbf{S}_i^T \mathcal{W}_i \mathbf{S}_i)^{-1} \mathbf{S}_i^T \mathcal{W}_i \boldsymbol{\eta}.$$

By the standard argument in Fan and Gijbels [3] and Lemma 1 in Fan and Zhang [4], we have

$$\mathcal{T} = (2K^2(0) - v_*^2)h^{-2} + o(h^{-2}),$$

when  $h = o(n^{-1/6})$  and  $nh^2 \rightarrow \infty$ , where  $v_* = \int K^2(t) dt$ .

We conclude that an unknown bivariate function amounts to  $(2K^2(0) - v_*^2)h^{-2}$  unknown constants. Based on this conclusion, if the number of constant components in  $\boldsymbol{\beta}(\cdot)$  is  $q$ , the  $\mathcal{K}$  in (3.1) will be  $q + (p - q)(2K^2(0) - v_*^2)h^{-2}$ .

To identify the constant components in  $\boldsymbol{\beta}(\cdot)$  in (1.3) is basically a model selection problem. Theoretically speaking, we go for the model with the smallest AIC (or BIC). However, in practice, it is almost computationally impossible to compute the AICs for all possible models. We have to use some algorithm to reduce the computational burden. In the following, we are going to introduce two algorithms for the model selection.

**3.2. Computational algorithms.** In this section, we use AIC as an example to demonstrate the introduced algorithms. The model in which  $\boldsymbol{\beta}(\cdot)$  has its  $i_1$ th,  $i_2$ th,  $\dots$ ,  $i_k$ th components being constant is denoted by  $\{i_1, \dots, i_k\}$ . When  $k = 0$ , we define the model as the model in which all components of  $\boldsymbol{\beta}(\cdot)$  are functional, and denote it by  $\{\}$ .

**Backward elimination.** The first algorithm we introduce is the backward elimination. Details are as follows.

(1) We start with the full model,  $\{1, \dots, p\}$ , and compute its AIC by (3.1). Denote the full model by  $\mathcal{M}_p$ , its AIC by  $\text{AIC}_p$ .

(2) For any integer  $k$ , suppose the current model is  $\mathcal{M}_k = \{i_1, \dots, i_k\}$  with AIC given by  $\text{AIC}_k$ . Take  $\mathcal{M}_{k-1}$  to be the model with the largest maximum of log likelihood function among the models  $\{i_1, \dots, i_{j-1}, i_{j+1}, \dots, i_k\}$ ,  $j = 1, \dots, k$ . If  $\text{AIC}_k < \text{AIC}_{k-1}$ , the chosen model is  $\mathcal{M}_k$ , and the model selection is ended; otherwise, continue to compute  $\mathcal{M}_l$  and  $\text{AIC}_l$  until either  $\text{AIC}_l < \text{AIC}_{l-1}$  or  $l = 0$ .

*Curvature-to-average ratio (CTAR) based method.* A more aggressive way to reduce the computational burden involved in the model selection procedure is based on the ratio of the curvature of the estimated function to its average. Explicitly, we first treat all  $\beta_j(\cdot)$ ,  $j = 1, \dots, p$ , as functional. For each  $j$ ,  $j = 1, \dots, p$ , we compute the curvature-to-average ratio (CTAR)  $R_j$  of the estimated function  $\hat{\beta}_j(\cdot)$ :

$$R_j = \frac{1}{\bar{\beta}_j^2} \sum_{i=1}^n \{\hat{\beta}_j(s_i) - \bar{\beta}_j\}^2, \quad \bar{\beta}_j = \frac{1}{n} \sum_{i=1}^n \hat{\beta}_j(s_i), \quad j = 1, \dots, p.$$

We sort  $R_j$ ,  $j = 1, \dots, p$ , in an increasing order, say  $R_{i_1} \leq \dots \leq R_{i_p}$ , then compute the AICs for the models  $\{i_1, \dots, i_k\}$  from  $k = 0$  to the turning point  $k_0$  where the AIC starts to increase. The chosen model is  $\{i_1, \dots, i_{k_0}\}$ .

The algorithm based on the CTAR is much faster than the backward elimination based algorithm, however, we find it less accurate although it still works reasonably well in our simulation studies. This is because the CTARs of all coefficients are obtained in one go based on the model in which all coefficients are treated as functional, and not updated. This will speed up the selection procedure; on the other hand, the effect of randomness would be stronger than that in backward elimination, which leads to a slightly larger possibility of picking up a wrong model.

**4. Asymptotic properties.** In this section, we are going to present the asymptotic properties of the proposed estimators. We will, in this section, only present the asymptotic results, and leave the theoretical proofs in the [Appendix](#).

Although we assume  $\varepsilon_i$  in (1.3) follows normal distribution in our model assumption, we do not need this assumption when deriving the asymptotic properties of the proposed estimators. So, in this section, we do not assume  $\varepsilon_i$  follows normal distribution unless otherwise stated.

In this section, for  $w_{ij}$  in (1.3), we assume that there exists a sequence  $\rho_n > 0$  such that  $w_{ij} = O(1/\rho_n)$  uniformly with respect to  $i, j$  and the matrices  $W$  and  $A^{-1}$  are uniformly bounded in both row and column sums.

We now introduce some notations needed in the presentation of the asymptotic properties of the proposed estimators: let  $\mu_j = E\varepsilon_1^j$ ,  $j = 1, \dots, 4$ ,

$$\begin{aligned} \kappa_0 &= \int_{R^2} K(\|s\|) ds, \\ \kappa_2 &= \int_{R^2} [(1, 0)s]^2 K(\|s\|) ds = \int_{R^2} [(0, 1)s]^2 K(\|s\|) ds, \end{aligned}$$



$$\begin{aligned}
v_0 &= \int_{R^2} K^2(\|s\|) ds, \\
v_2 &= \int_{R^2} [(1, 0)s]^2 K^2(\|s\|) ds = \int_{R^2} [(0, 1)s]^2 K^2(\|s\|) ds, \\
G &= (g_{ij}) = WA^{-1}, \quad \Psi = E(X_1 X_1^T), \quad \Gamma = EX_1, \\
Z_1(s) &= \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n g_{ii} \beta(s_i) K_h(\|s_i - s\|), \\
Z_2(s) &= \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n \sum_{j \neq i}^n g_{ij} \beta(s_j) K_h(\|s_i - s\|), \\
Z(s) &= Z_1(s) + \Psi^{-1} \Gamma \Gamma^T Z_2(s), \\
Z &= \kappa_0^{-1} (f^{-1}(s_1) X_1^T Z(s_1), \dots, f^{-1}(s_n) X_n^T Z(s_n))^T, \\
\pi_1 &= \lim_{n \rightarrow \infty} \frac{\text{tr}((G + G^T)G)}{n}, \quad \pi_2 = \lim_{n \rightarrow \infty} \frac{\text{tr}(G)}{n}, \\
\pi_3 &= \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n g_{ii}^2, \\
\lambda_1 &= \lim_{n \rightarrow \infty} \frac{1}{n} E[(G\mathbf{m} - Z)^T (G\mathbf{m} - Z)], \\
\lambda_2 &= \lim_{n \rightarrow \infty} \frac{1}{n} E[(G\mathbf{m} - Z)^T G_c], \quad \lambda_3 = \lim_{n \rightarrow \infty} \frac{1}{n} E[(G\mathbf{m} - Z)^T \mathbf{1}_n],
\end{aligned}$$

where  $G_c = (g_{11}, \dots, g_{nn})^T$  and  $\mathbf{1}_n$  is an  $n$ -dimensional vector with each component being 1. Further, let

$$\begin{aligned}
\Omega &= \begin{pmatrix} \frac{1}{\sigma^2} \lambda_1 + \pi_1 & \frac{1}{\sigma^2} \pi_2 \\ \frac{1}{\sigma^2} \pi_2 & \frac{1}{2\sigma^4} \end{pmatrix}, \\
\Sigma &= \begin{pmatrix} \frac{\mu_4 - 3\sigma^4}{\sigma^4} \pi_3 + \frac{2\mu_3}{\sigma^4} \lambda_2 & \frac{\mu_3}{2\sigma^6} \lambda_3 + \frac{\mu_4 - 3\sigma^4}{2\sigma^6} \pi_2 \\ \frac{\mu_3}{2\sigma^6} \lambda_3 + \frac{\mu_4 - 3\sigma^4}{2\sigma^6} \pi_2 & \frac{\mu_4 - 3\sigma^4}{4\sigma^8} \end{pmatrix}, \\
s &= (u, v)^T, \quad \beta_{uu}(s) = \left( \frac{\partial^2 \beta_1(s)}{\partial u^2}, \dots, \frac{\partial^2 \beta_p(s)}{\partial u^2} \right)^T, \\
\beta_{vv}(s) &= \left( \frac{\partial^2 \beta_1(s)}{\partial v^2}, \dots, \frac{\partial^2 \beta_p(s)}{\partial v^2} \right)^T
\end{aligned}$$

and

$$S = \begin{pmatrix} (X_1^T, \mathbf{0}_{1 \times 2p})(\mathcal{X}_{(1)}^T \mathcal{W}_{(1)} \mathcal{X}_{(1)})^{-1} \mathcal{X}_{(1)}^T \mathcal{W}_{(1)} \\ \vdots \\ (X_n^T, \mathbf{0}_{1 \times 2p})(\mathcal{X}_{(n)}^T \mathcal{W}_{(n)} \mathcal{X}_{(n)})^{-1} \mathcal{X}_{(n)}^T \mathcal{W}_{(n)} \end{pmatrix},$$

where  $\mathcal{X}_{(i)}$  and  $\mathcal{W}_{(i)}$  are  $\mathcal{X}$  and  $\mathcal{W}$ , respectively, with  $s$  being replaced by  $s_i$ ,  $i = 1, \dots, n$ .

By some simple calculations, we can see the matrix  $\Omega$  defined above is the limit of the Fisher information matrix of  $\alpha$  and  $\sigma^2$ . As the singularity of matrix  $\Omega$  may have serious implication on the convergence rate of the proposed estimators, we present the asymptotic properties for the case where  $\Omega$  is nonsingular and the case where  $\Omega$  is singular separately. We present the nonsingular case in Theorems 1–3, and singular case in Theorems 4–7.

**THEOREM 1.** *Under the conditions (1)–(7) or conditions (1)–(6), (7̃) and (8) in Appendix,  $\Omega$  is nonsingular, and when  $n^{1/2}h^2/\log^2 n \rightarrow \infty$  and  $nh^8 \rightarrow 0$ ,  $\hat{\alpha}$  and  $\hat{\sigma}^2$  are consistent estimators of  $\alpha$  and  $\sigma^2$ , respectively.*

Theorem 1 shows the conditions under which  $\Omega$  is nonsingular and the consistency of  $\hat{\alpha}$  and  $\hat{\sigma}^2$  under such conditions. Based on Theorem 1, we can derive the asymptotic normality of  $\hat{\alpha}$  and  $\hat{\sigma}^2$ .

**THEOREM 2.** *Under the assumptions of Theorem 1, if the second partial derivative of  $\beta(s)$  is Lipschitz continuous and  $nh^6 \rightarrow 0$ ,*

$$\sqrt{n}(\hat{\alpha} - \alpha, \hat{\sigma}^2 - \sigma^2)^T \xrightarrow{D} N(\mathbf{0}, \Omega^{-1} + \Omega^{-1} \Sigma \Omega^{-1}).$$

Further, if  $\varepsilon_i$  is normally distributed,

$$\sqrt{n}(\hat{\alpha} - \alpha, \hat{\sigma}^2 - \sigma^2)^T \xrightarrow{D} N(\mathbf{0}, \Omega^{-1}).$$

Theorem 2 implies that the convergence rate of  $\hat{\alpha}$  is of order  $n^{-1/2}$  when  $\Omega$  is nonsingular, which is the optimal rate for parametric estimation. We will see, in Theorem 5, this rate can not be achieved by  $\hat{\alpha}$  when  $\Omega$  is singular.

**THEOREM 3.** *Under the assumptions of Theorem 1, if  $nh_1^6 = O(1)$  and  $h/h_1 \rightarrow 0$ ,*

$$\begin{aligned} & \sqrt{nh_1^2 f(s)} (\hat{\beta}(s) - \beta(s) - 2^{-1} \kappa_0^{-1} \kappa_2 h_1^2 \{\beta_{uu}(s) + \beta_{vv}(s)\}) \\ & \xrightarrow{D} N(\mathbf{0}, \kappa_0^{-2} v_0 \sigma^2 \Psi^{-1}) \end{aligned}$$

for any given  $s$ .

Theorem 3 shows  $\hat{\beta}(\cdot)$  is asymptotic normal and achieves the convergence rate of order  $n^{-1/6}$ , which is the optimal rate for bivariate nonparametric estimation.

We now turn to the case where  $\Omega$  is singular.

**THEOREM 4.** *Under the conditions (1)–(6) and (9) in the [Appendix](#),  $\Omega$  is singular, and if  $nh^8 \rightarrow 0$ ,  $n^{1/2}h^2/\log^2 n \rightarrow \infty$ ,  $\rho_n \rightarrow \infty$ ,  $\rho_n h^4 \rightarrow 0$  and  $nh^2/\rho_n \rightarrow \infty$ ,  $\hat{\alpha}$  is a consistent estimator of  $\alpha$ .*

**THEOREM 5.** *Under the assumptions of Theorem 4, if the second partial derivative of  $\beta(s)$  is Lipschitz continuous and  $nh^6 \rightarrow 0$ ,*

$$\sqrt{n/\rho_n}(\hat{\alpha} - \alpha) \xrightarrow{D} N(0, \sigma^2 \lambda_4^{-1}),$$

where

$$\lambda_4 = \lim_{n \rightarrow \infty} \frac{\rho_n}{n} E[(G\mathbf{m} - SG\mathbf{m})^T (G\mathbf{m} - SG\mathbf{m})].$$

Theorem 5 shows the convergence rate of  $\hat{\alpha}$  is of order  $(n/\rho_n)^{-1/2}$  which is slower than  $n^{-1/2}$  when  $\rho_n \rightarrow \infty$ . However, we will see, from Theorem 7, this has no effect on the asymptotic properties of  $\hat{\beta}(\cdot)$ .

**THEOREM 6.** *Under the assumptions of Theorem 5,*

$$\sqrt{n}(\hat{\sigma}^2 - \sigma^2) \xrightarrow{D} N(0, \mu_4 - \sigma^4).$$

Theorem 6 shows that although the asymptotic variance of  $\hat{\sigma}^2$  is different to that when  $\Omega$  is nonsingular,  $\hat{\sigma}^2$  still enjoys convergence rate of  $n^{-1/2}$ .

**THEOREM 7.** *Under the assumptions of Theorem 4, if  $nh_1^6 = O(1)$  and  $h/h_1 \rightarrow 0$ ,*

$$\begin{aligned} & \sqrt{nh_1^2 f(s)} (\hat{\beta}(s) - \beta(s) - 2^{-1} \kappa_0^{-1} \kappa_2 h_1^2 \{\beta_{uu}(s) + \beta_{vv}(s)\}) \\ & \xrightarrow{D} N(\mathbf{0}, \kappa_0^{-2} v_0 \sigma^2 \Psi^{-1}) \end{aligned}$$

for any given  $s$ .

From Theorems 3 and 7, we can see the singularity of  $\Omega$  has no effect on the asymptotic distribution of  $\hat{\beta}(\cdot)$ .

**5. Simulation studies.** In this section, we will use simulated examples to examine the performances of the proposed estimation and model selection procedure. In all simulated examples and the real data analysis later on, we set  $w_{ij}$  to be

$$(5.1) \quad w_{ij} = \exp(-\|s_i - s_j\|) / \sum_{k \neq i} \exp(-\|s_i - s_k\|).$$

We first examine the performance of the proposed estimation procedure, then the model selection procedure.

### 5.1. Performance of the estimation procedure.

EXAMPLE 1. In model (1.3), we set  $p = 3$ ,  $\sigma^2 = 1$ ,

$$\alpha = 0.5, \quad \beta_1(s) = \sin(\|s\|^2\pi), \quad \beta_2(s) = \cos(\|s\|^2\pi), \quad \beta_3(s) = e^{(\|s\|^2)}$$

and independently generate  $X_i$  from  $N(\mathbf{0}_3, I_3)$ ,  $s_i$  from  $U[0, 1]^2$ ,  $\varepsilon_i$  from  $N(0, \sigma^2)$ ,  $i = 1, \dots, n$ .  $y_i$ ,  $i = 1, \dots, n$ , are generated through model (1.3). We are going to apply the proposed estimation method based on the generated  $(s_i, X_i^T, y_i)$ ,  $i = 1, \dots, n$ , to estimate  $\beta_1(\cdot)$ ,  $\beta_2(\cdot)$ ,  $\beta_3(\cdot)$ ,  $\alpha$  and  $\sigma^2$ , and examine the accuracy of the proposed estimation procedure.

We use the Epanechnikov kernel  $K(t) = 0.75(1 - t^2)_+$  as the kernel function in the estimation procedure. The bandwidth used in the estimation is 0.4.

We use mean squared error (MSE) to assess the accuracy of an estimator of an unknown constant parameter, mean integrated squared error (MISE) to assess the accuracy of an estimator of an unknown function.

For each given sample size  $n$ , we do 200 simulations. We compute the MSEs of the estimators of the unknown constants and the MISEs of the estimators of the unknown functions for sample size  $n = 400$ ,  $n = 500$  and  $n = 600$ . The obtained results are presented in Table 1. Table 1 shows the proposed estimation procedure works very well. To have a more visible idea about the performance of the proposed estimation procedure, we set sample size  $n = 500$  and do 200 simulations. We single out the one with median performance among the 200 simulations. The estimate of  $\alpha$  coming from this simulation is 0.407, the estimate of  $\sigma^2$  is 0.976. The estimated unknown functions from this simulation are presented in Figures 1, 2 and 3, and are superimposed with the true functions. All these show our estimation procedure works very well.

### 5.2. Performance of the model selection procedure.

EXAMPLE 2. In model (1.3), we set  $p = 5$ ,  $\beta_1(\cdot)$ ,  $\beta_2(\cdot)$  and  $\beta_3(\cdot)$  the same as that in Example 1,  $\beta_4(\cdot) = \sin^2(\|s\|^2\pi)$ ,  $\beta_5(\cdot) = \beta_5 = 1$ . We generate  $X_i$ ,  $s_i$ ,  $\varepsilon_i$ ,  $y_i$   $i = 1, \dots, n$ , in the same way as that in Example 1, except that  $X_i$  is from

TABLE 1  
The MISEs and MSEs

	$\hat{\beta}_1(\cdot)$	$\hat{\beta}_2(\cdot)$	$\hat{\beta}_3(\cdot)$	$\hat{\alpha}$	$\hat{\sigma}^2$
$n = 400$	0.0769	0.0642	0.0618	0.0128	0.0086
$n = 500$	0.0712	0.0573	0.0539	0.0093	0.0065
$n = 600$	0.0679	0.0498	0.0474	0.0076	0.0053

The column corresponding to the estimator of an unknown function is the MISEs of the estimator for  $n = 400$ ,  $n = 500$  and  $n = 600$ , corresponding to the estimator of an unknown constant is the MSEs of the estimator.

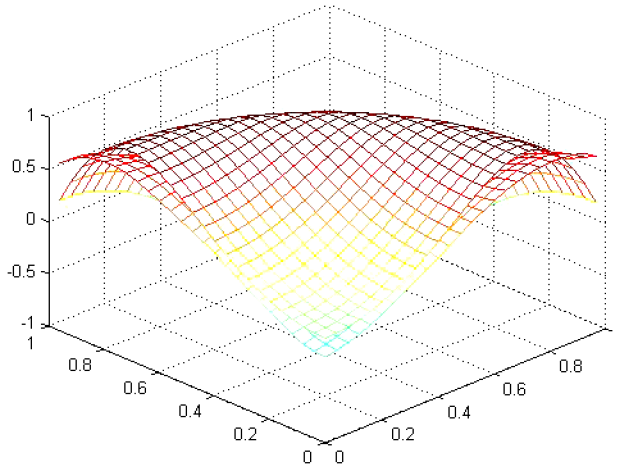


FIG. 1. The estimated  $\beta_1(s)$  superimposed with  $\beta_1(s)$ .

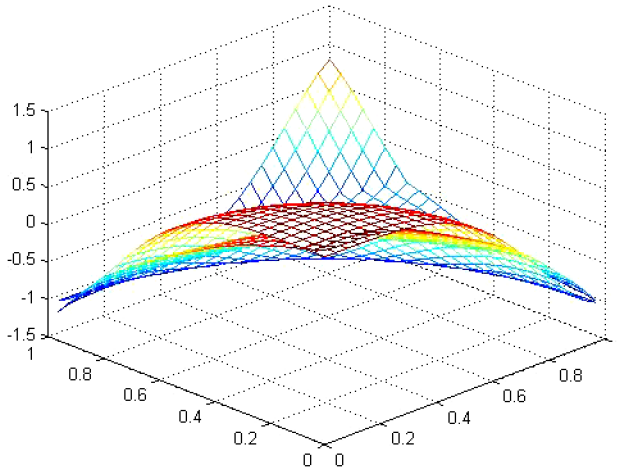


FIG. 2. The estimated  $\beta_2(s)$  superimposed with  $\beta_2(s)$ .

$N(\mathbf{0}_5, I_5)$ . Based on the generated data, we are going to apply the proposed AIC or BIC to select the correct model, and examine the performances of the proposed AIC, BIC and the two algorithms in identifying the constant components in model (1.3).

We still use the Epanechnikov kernel as the kernel function in the model selection, however, the bandwidth used is 0.2 for AIC and 0.3 for BIC, which is smaller than that for estimation. In general, the bandwidth used for model selection should be smaller than that for estimation. In fact, we have tried different bandwidths, it turned out any bandwidth in a reasonable range such as  $[0.15, 0.3]$  for AIC,  $[0.2, 0.35]$  for BIC would do the job very well.

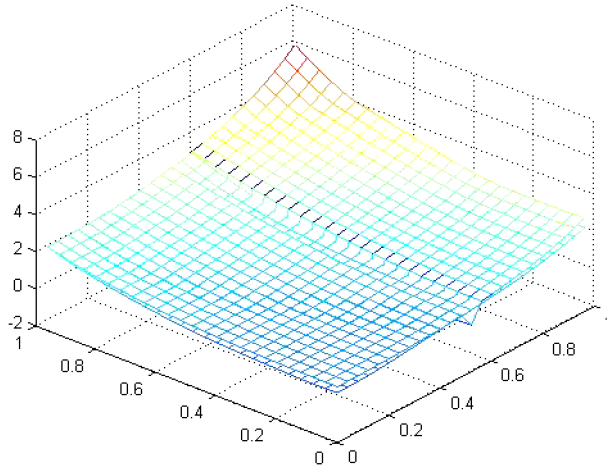


FIG. 3. The estimated  $\hat{\beta}_3(s)$  superimposed with  $\beta_3(s)$ .

Due to the very expensive computation involved, for any given sample size  $n$ , we only do 200 simulations, and in each simulation, we apply either AIC or BIC coupled with either of the two proposed algorithms to select model. For each candidate model, the ratios of picking up this model in the 200 simulations are computed for different cases. The results are presented in Table 2. We can see, from Table 2, the proposed BIC with backward elimination performs best, and the others are doing reasonably well, also.

**6. Real data analysis.** In this section, we are going to apply the proposed model (1.3) together with the proposed model selection and estimation method to analyze the Boston house price data. Specifically, we are going to explore how some factors such as the per capita crime rate by town (denoted by CRIM), average number of rooms per dwelling (denoted by RM), index of accessibility to radial highways (denoted by RAD), full-value property-tax rate per \$10,000 dollar (denoted by TAX) and the percentage of the lower status of the population (denoted by LSTAT) affect the median value of owner-occupied homes in \$1000's (denoted by MEDV), and whether the effects of these factors vary over location or not.

We use model (1.3) to fit the data with  $y_i$ ,  $x_{i1}$ ,  $x_{i2}$ ,  $x_{i3}$ ,  $x_{i4}$  and  $x_{i5}$  being MEDV, CRIM, RM, RAD, TAX and LSTAT, respectively, and  $X_i = (x_{i1}, \dots, x_{i5})^T$ . The kernel function used in either estimation procedure or model selection is taken to be the Epanechnikov kernel.

We first try to find which factors have location varying effects on the house price, and which factors do not. This is equivalent to identifying the constant coefficients in the model used to fit the data. We apply the proposed BIC coupled with backward elimination to do the model selection, and the bandwidth used is chosen to be 17% of the range of the locations. The obtained result shows the coefficients

TABLE 2  
Ratios of picking up each model in model selection

	{5}	{1, 5}	{4, 5}	{1, 4, 5}	{1, 2, 4, 5}	{1, 2, 3, 4, 5}
$n = 400$	0.83	0.05	0.07	0.02	0.02	0.01
$n = 500$	0.91	0.02	0.05	0.02	0	0
$n = 600$	0.94	0.02	0.03	0	0	0.01
$n = 400$	0.81	0.06	0.08	0.05	0	0
$n = 500$	0.89	0.03	0.06	0.02	0	0
$n = 600$	0.92	0.01	0.04	0.01	0.02	0
$n = 400$	0.86	0.04	0.05	0.03	0.01	0.01
$n = 500$	0.93	0.02	0.03	0.01	0.01	0
$n = 600$	0.96	0.01	0.02	0.01	0	0
$n = 400$	0.84	0.05	0.07	0.02	0.01	0.01
$n = 500$	0.88	0.05	0.05	0.01	0.01	0
$n = 600$	0.93	0.03	0.03	0.01	0	0

The ratios of picking up each candidate model in 200 simulations for different sample sizes.  $\{i_1, \dots, i_k\}$  stands for the model in which  $\beta(\cdot)$  has its  $i_1$ th,  $\dots$ ,  $i_k$ th components being constant and the column corresponding to which is the ratios of picking up this model among 200 simulations. Row 2 to row 4 are the ratios obtained based on AIC and backward elimination when sample size  $n = 400$ ,  $n = 500$  and  $n = 600$ . Row 5 to row 7 are the ratios obtained based on AIC and the CTAR based algorithm, row 8 to row 10 are the ratios obtained based on BIC and backward elimination, and row 11 to row 13 are the ratios obtained based on BIC and the CTAR based algorithm.

of  $x_{i3}$  and  $x_{i5}$  are constant, which means all factors, except RAD and LSTAT, have location varying effects on the house price.

We now apply the chosen model

$$(6.1) \quad y_i = \alpha \sum_{j \neq i} w_{ij} y_j + x_{i1} \beta_1(s_i) + x_{i2} \beta_2(s_i) + x_{i3} \beta_3 + x_{i4} \beta_4(s_i) + x_{i5} \beta_5 + \varepsilon_i,$$

$i = 1, \dots, n$ , where  $w_{ij}$  is defined by (5.1), to fit the data. The sample size of this data set is  $n = 506$ . The proposed estimation procedure is used to estimate the unknown functions and constants, and the bandwidth used in the estimation procedure is taken to be 60% of the range of the locations. The estimates of the unknown constants are presented in Table 3, and the estimates of the unknown functions are presented in Figure 4.

TABLE 3  
Estimates of the unknown constant coefficients

$\hat{\alpha}$	$\hat{\beta}_3$	$\hat{\beta}_5$
0.2210	0.3589	-0.4473

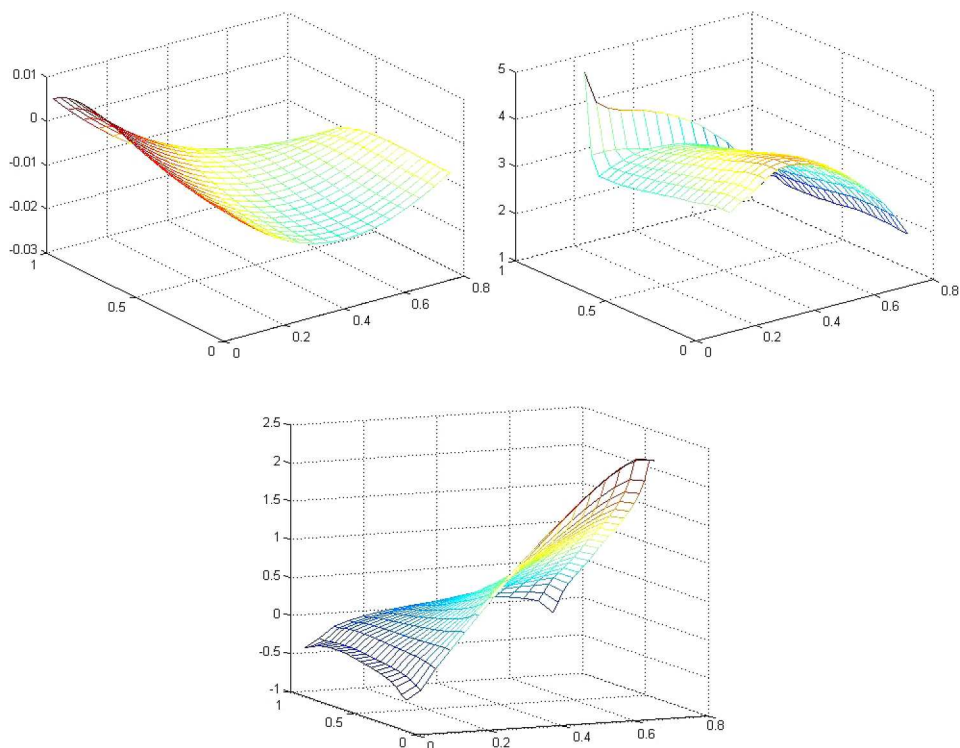


FIG. 4. The 3D plots of  $\hat{\beta}_1(s)$ ,  $\hat{\beta}_2(s)$  and  $\hat{\beta}_4(s)$ . The left one in the upper panel is  $\hat{\beta}_1(s)$ , right one in the upper panel is  $\hat{\beta}_2(s)$ , and the one in the lower panel is  $\hat{\beta}_4(s)$ .

To see how well model (6.1) fits the data, we conduct some residual analysis. The plot, normal  $Q-Q$  plot, ACF and partial ACF of the residuals of the fitting are presented in Figure 5. Figure 5 shows model (6.1) fits the data well.

As  $\beta_3$  and  $\beta_5$  can be interpreted as the impacts of RAD and LSTAT, respectively, Table 3 shows the index of accessibility to radial highways has positive impact on house price and the percentage of the lower status of the population has negative impact on house price. Apparently, this makes sense. Table 3 also shows that the estimate of  $\alpha$  is 0.221, which is an unignorable effect, and indicates the house prices in a neighborhood do affect each other. This is a true phenomenon in real world.

From Figure 4, we can see the impact  $\beta_1(\cdot)$  of the per capita crime rate by town on house price is negative and is clearly varying over location. The impact  $\beta_2(\cdot)$  of the average number of rooms per dwelling on house price is positive and is also varying over location. It is interesting to see that the impact of the average number of rooms per dwelling is lower in the area where the impact of crime rate is high than the area where the impact of crime rate is low. This implies that the crime rate is a dominate factor on the house price in the area where the impact of crime



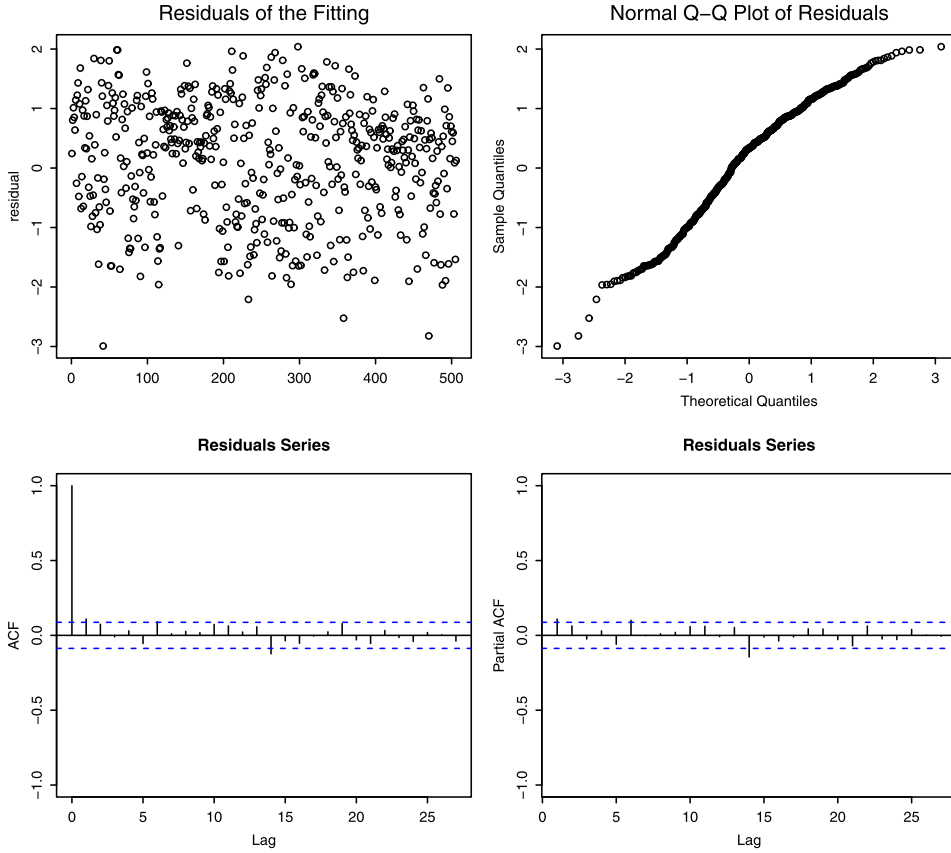


FIG. 5. The plot, normal  $Q$ - $Q$  plot, ACF and partial ACF of the residuals of the fitting of (6.1) to the Boston house price data.

rate is high. Figure 4 also shows the association between the house price and the full-value property-tax rate is varying over location, and it is generally positive, however, there are some areas where this association is negative. We can also see that the impact of the average number of rooms per dwelling is lower in the area, where the association between the house price and the full-value property tax rate is strong, than the area where the association is weak.

#### APPENDIX: CONDITIONS AND SKETCH OF THEORETICAL PROOFS

To avoid confusion of notation, we use  $\alpha_0$  to denote the true value of  $\alpha$  in this section. Further, we rewrite  $A = I_n - \alpha W$  as  $A(\alpha)$  to emphasis its dependence on  $\alpha$  and abbreviate  $A(\alpha_0)$  as  $A$ .

The following regularity conditions are needed to establish the asymptotic properties of the estimators.

**Conditions.**

(1) The kernel function  $K(\cdot)$  is a bounded positive, symmetric and Lipschitz continuous function with a compact support on  $\mathbb{R}$ .  $h \rightarrow 0$ .

(2)  $\{\beta_i(\cdot), i = 1, \dots, p\}$  have continuous second partial derivatives.

(3)  $\{X_1, \dots, X_n\}$  is an i.i.d. random sample and is independent of  $\{\varepsilon_1, \dots, \varepsilon_n\}$ . Moreover,  $E(X_1 X_1^T)$  is positive definite,  $E\|X_1\|^{2q} < \infty$  and  $E|\varepsilon_1|^{2q} < \infty$  for some  $q > 2$ .

(4)  $\{s_i\}$  is a sequence of fixed design points on a bounded compact support  $\mathcal{S}$ . Further, there exists a positive joint density function  $f(\cdot)$  satisfying a Lipschitz condition such that

$$\sup_{s \in \mathcal{S}} \left| \frac{1}{n} \sum_{i=1}^n [r(s_i) K_h(\|s_i - s\|)] - \int r(t) K_h(\|t - s\|) f(t) dt \right| = O(h)$$

for any bounded continuous function  $r(\cdot)$  and  $K_h(\cdot) = K(\cdot/h)/h^2$  where  $K(\cdot)$  satisfies condition (1).  $f(\cdot)$  is bounded away from zero on  $\mathcal{S}$ .

(5)  $w_{ii} = 0$  for any  $i$ , and there exists a sequence  $\rho_n > 0$  such that  $w_{ij} = O(1/\rho_n)$  uniformly with respect to  $i$  and  $j$ . Furthermore, the matrices  $W$  and  $A^{-1}$  are uniformly bounded in both row and column sums.

(6)  $A^{-1}(\alpha)$  are uniformly bounded in either row or column sums, uniformly in  $\alpha$  in a compact support  $\Delta$ . The true  $\alpha_0$  is an interior point in  $\Delta$ .

(7)  $\lim_{n \rightarrow \infty} \frac{1}{n} E[(G\mathbf{m} - Z)^T(G\mathbf{m} - Z)] = \lambda_1 > 0$ .

(7)  $\lambda_1 = 0$ .

(8)  $\rho_n$  is bounded and for any  $\alpha \neq \alpha_0$ ,

$$\lim_{n \rightarrow \infty} \left\{ \frac{1}{n} \log |\sigma^2 A^{-1} (A^{-1})^T| - \frac{1}{n} \log |\sigma_a^2(\alpha) A^{-1}(\alpha) (A^{-1}(\alpha))^T| \right\} \neq 0,$$

where  $\sigma_a^2(\alpha) = \frac{\sigma^2}{n} \text{tr}\{(A(\alpha)A^{-1})^T A(\alpha)A^{-1}\}$ .

(9)  $\rho_n \rightarrow \infty$ , the row sums of  $G$  have the uniform order  $O(1/\sqrt{\rho_n})$  and

$$\lim_{n \rightarrow \infty} \frac{\rho_n}{n} E[(G\mathbf{m} - SG\mathbf{m})^T(G\mathbf{m} - SG\mathbf{m})] = \lambda_4 > 0.$$

**REMARK 1.** Conditions (1)–(3) are commonly seen in nonparametric estimation. They are not the weakest possible ones, but they are imposed to facilitate the technical proofs. Since the sampling units can be regarded as given, the fixed bounded design condition (4) is made for technical convenience. Of course, as in Linton [11], condition (4) does not preclude  $\{s_i\}_{i=1}^n$  from being generated by some random mechanism. For example, if  $s_i$ 's were i.i.d. with joint density  $f(\cdot)$ , then condition (4) holds with probability one which can be obtained in a similar way to Hansen [7]. So, we can obtain our results by firstly conditional on  $\{s_i\}_{i=1}^n$ , then some standard arguments.

REMARK 2. Conditions (5)–(8) parallel the corresponding conditions of Lee [9] and Su and Jin [13]. Conditions (5)–(6) concern the essential features of the weight matrix for the model. Condition (7) is a sufficient condition which ensures that the likelihood function of  $\alpha$  has a unique maximizer. When condition (7) holds and the elements of  $W$  are uniformly bounded, the uniqueness of the maximizer can be guaranteed by condition (8). These two kinds of conditions ensure that  $\Omega$  which is the limit of the information matrix of the finite-dimensional parameters is nonsingular. So, they are the crucial conditions for  $\sqrt{n}$ -rate of convergence of the finite-dimensional parameter estimators.

REMARK 3. When  $\rho_n \rightarrow \infty$ ,  $\Omega$  is nonsingular only when condition (7) holds. Under condition (7),  $\Omega$  will become singular. The singularity of the matrix may have implications on the rate of convergence of the estimators. Nevertheless, we follow Lee [9] and Su and Jin [13] to consider the situation where

$$\lim_{n \rightarrow \infty} \frac{\rho_n}{n} E[(G\mathbf{m})^T (I_n - S)^T (I_n - S) G\mathbf{m}] = \lambda_4 \in (0, \infty).$$

In this case, it is natural to assume that the elements of  $(I_n - S)G\mathbf{m}$  have the uniform order  $O_P(1/\sqrt{\rho_n})$  which can be satisfied by the assumption that the row sums of  $G$  are of uniform order  $O(1/\sqrt{\rho_n})$ .

In the following, let  $H$  be a diagonal matrix of size  $3p$  with its first  $p$  elements on the diagonal being 1 and the remaining elements being  $h$ ,  $P = (I_n - S)^T (I_n - S)$  and  $\boldsymbol{\varepsilon} = (\varepsilon_1, \dots, \varepsilon_n)^T$ . Moreover, like  $\alpha_0$ , we use  $\sigma_0^2$  to denote the true value of  $\sigma^2$  to avoid confusion of notation. Since the following notations will be frequently used in the proofs, we list here for easy reference:

$$\begin{aligned} l(\alpha, \sigma^2) &= -\frac{n}{2} \log(\sigma^2) + \log(|A(\alpha)|) - \frac{1}{2\sigma^2} (A(\alpha)Y)^T P A(\alpha)Y, \\ l_c(\alpha) &= -\frac{n}{2} \log \tilde{\sigma}^2(\alpha) + \log |A(\alpha)|, \\ \tilde{\sigma}^2(\alpha) &= \frac{1}{n} (A(\alpha)Y)^T P A(\alpha)Y, \\ \bar{\sigma}^2(\alpha) &= \frac{1}{n} E[(A(\alpha)Y)^T P A(\alpha)Y], \\ \sigma_a^2(\alpha) &= \frac{\sigma_0^2}{n} \text{tr}\{(A(\alpha)A^{-1})^T A(\alpha)A^{-1}\}. \end{aligned}$$

To prove the theorems, the following lemmas are needed. Their proofs and the more detailed proofs of the theorems can be found in the supplementary material (Sun et al. [14]).

LEMMA 1. Let  $\{Y_i\}$  be a sequence of independent random variables and  $\{s_i\} \in \mathbb{R}^2$  are nonrandom vectors. Suppose that for some  $q > 2$ ,  $\max_i E|Y_i|^q < \infty$ . Then under condition (1), we have

$$\sup_{s \in \mathcal{S}} \left| \frac{1}{n} \sum_{i=1}^n [K_h(\|s_i - s\|) Y_i - E\{K_h(\|s_i - s\|) Y_i\}] \right| = O_p\left(\left\{\frac{\log n}{nh^2}\right\}^{1/2}\right),$$

provided that  $n^{1-2/q}h^2/\log^2 n \rightarrow \infty$  and  $\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n K_h(\|s_i - s\|) < \infty$  for any  $s \in \mathcal{S}$ .

LEMMA 2. Under conditions (1)–(4), when  $n^{1/2}h^2/\log^2 n \rightarrow \infty$ ,

$$(1) \quad n^{-1}H^{-1}\mathcal{X}^T\mathcal{W}\mathcal{X}H^{-1} = \begin{pmatrix} \kappa_0 f(s)\Psi & \mathbf{0}_{p \times 2p} \\ \mathbf{0}_{2p \times p} & \kappa_2 f(s)\Psi \otimes I_2 \end{pmatrix} + O_p(c_n \mathbf{1}_{3p} \mathbf{1}_{3p}^T)$$

holds uniformly in  $s \in \mathcal{S}$  where  $c_n = h + \{\frac{\log n}{nh^2}\}^{1/2}$ ,

$$(2) \quad \boldsymbol{\beta}(s) - (I_p, \mathbf{0}_{p \times 2p})(\mathcal{X}^T\mathcal{W}\mathcal{X})^{-1}\mathcal{X}^T\mathcal{W}\mathbf{m} = -\frac{\kappa_2 h^2}{2\kappa_0} \{\boldsymbol{\beta}_{uu}(s) + \boldsymbol{\beta}_{vv}(s)\} + o_p(h^2 \mathbf{1}_p)$$

holds uniformly in  $s \in \mathcal{S}$ .

LEMMA 3. Under conditions (1)–(5), when  $n^{1/2}h^2/\log^2 n \rightarrow \infty$ ,

$$n^{-1}H^{-1}\mathcal{X}^T\mathcal{W}\mathbf{G}\mathbf{m} - n^{-1}E(H^{-1}\mathcal{X}^T\mathcal{W}\mathbf{G}\mathbf{m}) = o_p(1)$$

uniformly in  $s \in \mathcal{S}$ .

LEMMA 4. Under conditions (1), (3), (4) and (5), when  $n^{1/2}h^2/\log^2 n \rightarrow \infty$ , we have (1)  $\frac{1}{n}E[\text{tr}(P)] = 1 + o(1)$ , (2)  $\frac{1}{n}E[\text{tr}(G^T P) - \text{tr}(G)] = o(1)$ , (3)  $\frac{1}{n}E[\text{tr}(G^T P G) - \text{tr}(G^T G)] = o(1)$ . Further, when  $nh^2/\rho_n \rightarrow \infty$ , (4)  $\frac{\rho_n}{n}E[\text{tr}(P) - n] = o(1)$ , (5)  $\frac{\rho_n}{n}E[\text{tr}(G^T P) - \text{tr}(G)] = o(1)$ , (6)  $\frac{\rho_n}{n}E[\text{tr}(G^T P G) - \text{tr}(G^T G)] = o(1)$ .

LEMMA 5. Under conditions (1)–(5), when  $n^{1/2}h^2/\log^2 n \rightarrow \infty$ , (1)  $(\mathbf{G}\mathbf{m})^T P \mathbf{m} = o_p(nh^2)$ . Moreover, under the assumption that the second partial derivative of  $\boldsymbol{\beta}(s)$  is Lipschitz continuous, we have (2)  $(\mathbf{G}\mathbf{m})^T P \mathbf{m} = O_p(nh^3 + \{nh^2 \log n\}^{1/2})$ .

LEMMA 6. Under conditions (1)–(5), when  $n^{1/2}h^2/\log^2 n \rightarrow \infty$  and  $nh^8 \rightarrow 0$ , we have (1)  $n^{-1/2}L^T P \mathbf{m} = o_p(1)$  for  $L = \mathbf{m}, \boldsymbol{\varepsilon}$  and  $G\boldsymbol{\varepsilon}$ , (2)  $n^{-1}L^T P G \mathbf{m} = o_p(1)$  for  $L = \mathbf{m}, \boldsymbol{\varepsilon}$  and  $G\boldsymbol{\varepsilon}$ .

LEMMA 7. Under conditions (1)–(5), when  $n^{1/2}h^2/\log^2 n \rightarrow \infty$ , we have (1)  $\frac{1}{n}\{(\mathbf{G}\mathbf{m})^T P \mathbf{G}\mathbf{m} - E[(\mathbf{G}\mathbf{m})^T P \mathbf{G}\mathbf{m}]\} = o_P(1)$ , (2)  $\frac{1}{n}E[(\mathbf{G}\mathbf{m})^T P \mathbf{G}\mathbf{m}] = \frac{1}{n}E[(\mathbf{G}\mathbf{m} - \mathbf{Z})^T(\mathbf{G}\mathbf{m} - \mathbf{Z})] + o(1)$ .

LEMMA 8. Under conditions (1)–(5), when  $n^{1/2}h^2/\log^2 n \rightarrow \infty$ , we have (1)  $n^{-1/2}\{\boldsymbol{\varepsilon}^T P \boldsymbol{\varepsilon} - \boldsymbol{\varepsilon}^T \boldsymbol{\varepsilon}\} = o_P(1)$ , (2)  $n^{-1/2}\{\boldsymbol{\varepsilon}^T G^T P \boldsymbol{\varepsilon} - \boldsymbol{\varepsilon}^T G^T \boldsymbol{\varepsilon}\} = o_P(1)$ , (3)  $n^{-1/2}\{\boldsymbol{\varepsilon}^T G^T P G \boldsymbol{\varepsilon} - \boldsymbol{\varepsilon}^T G^T G \boldsymbol{\varepsilon}\} = o_P(1)$ , (4)  $n^{-1/2}\{(\mathbf{G}\mathbf{m})^T P \boldsymbol{\varepsilon} - (\mathbf{G}\mathbf{m} - S\mathbf{G}\mathbf{m})^T \boldsymbol{\varepsilon}\} = o_P(1)$ .

LEMMA 9. Suppose that  $B = (b_{ij})_{1 \leq i, j \leq n}$  is a sequence of symmetric matrices with row and column sums uniformly bounded and its elements are also uniformly bounded. Let  $\sigma_{Q_n}^2$  be the variance of  $Q_n$  where  $Q_n = (\mathbf{G}\mathbf{m} - S\mathbf{G}\mathbf{m})^T \boldsymbol{\varepsilon} + \boldsymbol{\varepsilon}^T B \boldsymbol{\varepsilon} - \sigma_0^2 \text{tr}(B)$ . Assume that the variance  $\sigma_{Q_n}^2$  is  $O(n)$  with  $\{\frac{\sigma_{Q_n}^2}{n}\}$  bounded away from zero, then we have under conditions (1)–(5) that  $\frac{Q_n}{\sigma_{Q_n}} \xrightarrow{D} N(0, 1)$ .

LEMMA 10. Under conditions (1)–(5), and the row sums of matrix  $G$  having the uniform order  $O(1/\sqrt{\rho_n})$  and  $n^{1/2}h^2/\log^2 n \rightarrow \infty$ , we have (1)  $(\mathbf{G}\mathbf{m})^T P \mathbf{m} = o_P(\rho_n^{-1/2}nh^2)$ . Moreover, if the second partial derivative of  $\boldsymbol{\beta}(s)$  is Lipschitz continuous, then (2)  $(\mathbf{G}\mathbf{m})^T P \mathbf{m} = O_P(\rho_n^{-1/2}nh^3 + \{nh^2 \log n / \rho_n\}^{1/2})$ .

LEMMA 11. Under conditions (1)–(5) and the row sums of matrix  $G$  having the uniform order  $O(1/\sqrt{\rho_n})$ , when  $n^{1/2}h^2/\log^2 n \rightarrow \infty$ ,  $\rho_n \rightarrow \infty$ ,  $\rho_n h^4 \rightarrow 0$  and  $nh^2/\rho_n \rightarrow \infty$ , we have (1)  $\frac{\rho_n}{n} \mathbf{m}^T P \mathbf{m} = o_P(1)$ , (2)  $\frac{\rho_n}{n} L^T P \mathbf{G}\mathbf{m} = o_P(1)$  for  $L = \mathbf{m}, \boldsymbol{\varepsilon}$  and  $G\boldsymbol{\varepsilon}$ , (3)  $\sqrt{\frac{\rho_n}{n}}(G\boldsymbol{\varepsilon})^T P \mathbf{m} = o_P(1)$ , (4)  $\frac{\rho_n}{n}\{(\mathbf{G}\mathbf{m})^T P \mathbf{G}\mathbf{m} - E[(\mathbf{G}\mathbf{m})^T P \mathbf{G}\mathbf{m}]\} = o_P(1)$ , (5)  $\sqrt{\frac{\rho_n}{n}}\{\boldsymbol{\varepsilon}^T G^T P \boldsymbol{\varepsilon} - \boldsymbol{\varepsilon}^T G^T \boldsymbol{\varepsilon}\} = o_P(1)$ , (6)  $\sqrt{\frac{\rho_n}{n}}\{\boldsymbol{\varepsilon}^T G^T P G \boldsymbol{\varepsilon} - \boldsymbol{\varepsilon}^T G^T G \boldsymbol{\varepsilon}\} = o_P(1)$ , (7)  $\sqrt{\frac{\rho_n}{n}}\{(\mathbf{G}\mathbf{m})^T P \boldsymbol{\varepsilon} - (\mathbf{G}\mathbf{m} - S\mathbf{G}\mathbf{m})^T \boldsymbol{\varepsilon}\} = o_P(1)$ .

LEMMA 12. Suppose that  $B = (b_{ij})_{1 \leq i, j \leq n}$  is a sequence of symmetric matrices with row and column sums uniformly bounded. Let  $\sigma_{Q_n}^2$  be the variance of  $Q_n$  where  $Q_n = (\mathbf{G}\mathbf{m} - S\mathbf{G}\mathbf{m})^T \boldsymbol{\varepsilon} + \boldsymbol{\varepsilon}^T B \boldsymbol{\varepsilon} - \sigma_0^2 \text{tr}(B)$ . Assume that the variance  $\sigma_{Q_n}^2$  is  $O(n/\rho_n)$  with  $\{\frac{\rho_n}{n}\sigma_{Q_n}^2\}$  bounded away from zero, the elements of  $B$  are of uniform order  $O(1/\rho_n)$  and the row sums of  $G$  of uniform order  $O(1/\sqrt{\rho_n})$ , we have under  $\rho_n \rightarrow \infty$  and conditions (1)–(5) that  $\frac{Q_n}{\sigma_{Q_n}} \xrightarrow{D} N(0, 1)$ .

In the proofs of the theorems, we will use the facts that for constant matrices  $B = (b_{ij})$  and  $D = (d_{ij})$ ,  $\text{var}(\boldsymbol{\varepsilon}^T B \boldsymbol{\varepsilon}) = (\mu_4 - 3\sigma_0^4) \sum_{i=1}^n b_{ii}^2 + \sigma_0^4 [\text{tr}(B B^T) +$

$\text{tr}(B^2)]$  and

$$E(\boldsymbol{\varepsilon}^T B \boldsymbol{\varepsilon} \boldsymbol{\varepsilon}^T D \boldsymbol{\varepsilon}) = (\mu_4 - 3\sigma_0^4) \sum_{i=1}^n b_{ii} d_{ii} + \sigma_0^4 [\text{tr}(B) \text{tr}(D) + \text{tr}(BD) + \text{tr}(BD^T)].$$

Moreover, we will frequently use the following facts by condition (5) (see Lee [9]) without being clearly pointed out:

- (1) the elements of  $G = WA^{-1}$  are  $O(1/\rho_n)$  uniformly with respect to  $i$  and  $j$ .
- (2) The matrix  $G = WA^{-1}$  is uniformly bounded in both row and column sums.

**PROOF OF THEOREM 1.** We will first show that  $\Omega$  is nonsingular. Let  $\mathbf{d} = (d_1, d_2)^T$  be a constant vector such that  $\Omega \mathbf{d} = \mathbf{0}_2$ . Then it is sufficient to show that  $\mathbf{d} = \mathbf{0}_2$ . From the second equation of  $\Omega \mathbf{d} = \mathbf{0}_2$ , we have that  $d_2 = -2\sigma_0^2 \lim_{n \rightarrow \infty} \frac{1}{n} \text{tr}(G)d_1$ . Plugging  $d_2$  into the first equation of  $\Omega \mathbf{d} = \mathbf{0}_2$ , we have that

$$d_1 \left\{ \frac{1}{\sigma_0^2} \lambda_1 + \lim_{n \rightarrow \infty} \left[ \frac{1}{n} \text{tr}((G + G^T)G) - \frac{2}{n^2} \text{tr}^2(G) \right] \right\} = 0.$$

It follows by condition (7) that  $\lambda_1 > 0$ . Moreover,  $\text{tr}\{(G + G^T)G\} - \frac{2}{n} \text{tr}^2(G) = \frac{1}{2} \text{tr}\{(\tilde{G}^T + \tilde{G})(\tilde{G}^T + \tilde{G})^T\} \geq 0$  where  $\tilde{G} = G - \frac{1}{n} \text{tr}(G)I_n$ . As we have by condition (5) that  $\text{tr}\{(\tilde{G}^T + \tilde{G})(\tilde{G}^T + \tilde{G})^T\} = O(\frac{n}{\rho_n})$ , if condition (7) holds, condition (8) implies that the limit of  $\frac{1}{2n} \text{tr}\{(\tilde{G}^T + \tilde{G})(\tilde{G}^T + \tilde{G})^T\} > 0$ . Therefore,  $d_1 = 0$  and  $d_2 = 0$ .

Next, we will follow the idea of Lee [9] to show the consistency of  $\hat{\alpha}$ . Define  $Q(\alpha)$  to be  $\max_{\sigma^2} E[l(\alpha, \sigma^2)]$  by ignoring the constant term. The optimal solution of this maximization problem is  $\bar{\sigma}^2(\alpha) = \frac{1}{n} E[(A(\alpha)Y)^T P A(\alpha)Y]$ . Consequently,

$$Q(\alpha) = -n/2 \cdot \log \bar{\sigma}^2(\alpha) + \log |A(\alpha)|.$$

According to White ([18], Theorem 3.4), it suffices to show the uniform convergence of  $n^{-1}\{l_c(\alpha) - Q(\alpha)\}$  to zero in probability on  $\Delta$  and the unique maximizer condition that

$$(A.1) \quad \limsup_{n \rightarrow \infty} \max_{\alpha \in N^c(\alpha_0, \delta)} n^{-1}[Q(\alpha) - Q(\alpha_0)] < 0 \quad \text{for any } \delta > 0,$$

where  $N^c(\alpha_0, \delta)$  is the complement of an open neighborhood of  $\alpha_0$  in  $\Delta$  with diameter  $\delta$ .

Note that  $\frac{1}{n} l_c(\alpha) - \frac{1}{n} Q(\alpha) = -\frac{1}{2} \{\log \tilde{\sigma}^2(\alpha) - \log \bar{\sigma}^2(\alpha)\}$ , then to show the uniform convergence, it is sufficient to show that  $\tilde{\sigma}^2(\alpha) - \bar{\sigma}^2(\alpha) = o_P(1)$  uniformly on  $\Delta$  and  $\bar{\sigma}^2(\alpha)$  is uniformly bounded away from zero on  $\Delta$ .

As  $A(\alpha)A^{-1} = I_n + (\alpha_0 - \alpha)G$  by  $WA^{-1} = G$ , the result  $\tilde{\sigma}^2(\alpha) - \bar{\sigma}^2(\alpha) = o_P(1)$  uniformly on  $\Delta$  can be obtained by straightforward calculations, Lemmas 4(1)–(3), 6, 7(1), 8(1)–(3) and Chebyshev's inequality.

Now we will show that  $\bar{\sigma}^2(\alpha)$  is bounded away from zero uniformly on  $\Delta$ . As we know by simple calculations and Lemma 4(1)–(3) that

$$(A.2) \quad \bar{\sigma}^2(\alpha) \geq \sigma_0^2 n^{-1} \text{tr}\{(A(\alpha)A^{-1})^T A(\alpha)A^{-1}\} + o(1),$$

it suffices to show that  $\sigma_a^2(\alpha) = \frac{\sigma_0^2}{n} \text{tr}\{(A(\alpha)A^{-1})^T A(\alpha)A^{-1}\}$  is uniformly bounded away from zero on  $\Delta$ . To do so, we define an auxiliary spatial autoregressive (SAR) process:  $Y = \alpha_0 WY + \epsilon$  with  $\epsilon \sim N(\mathbf{0}, \sigma_0^2 I_n)$ . Its log likelihood function without the constant term is

$$l_a(\alpha, \sigma^2) = -\frac{n}{2} \log \sigma^2 + \log |A(\alpha)| - \frac{1}{2\sigma^2} (A(\alpha)Y)^T A(\alpha)Y.$$

Set  $Q_a(\alpha)$  to be  $\max_{\sigma^2} E_a[l_a(\alpha, \sigma^2)]$  by ignoring the constant term, where  $E_a$  is the expectation under this SAR process. It can be easily shown that

$$Q_a(\alpha) = -n/2 \cdot \log \sigma_a^2(\alpha) + \log |A(\alpha)|.$$

So, we have by Jensen's inequality that  $Q_a(\alpha) \leq Q_a(\alpha_0)$  for all  $\alpha \in \Delta$ , hence it follows:

$$-\frac{1}{2} \log \sigma_a^2(\alpha) \leq -\frac{1}{2} \log \sigma_0^2 + \frac{1}{n} (\log |A(\alpha_0)| - \log |A(\alpha)|)$$

uniformly on  $\Delta$ . Since we have, by the mean value theorem and conditions (5)–(6), that  $n^{-1} \{\log |A(\alpha_2)| - \log |A(\alpha_1)|\} = O(1)$  uniformly in  $\alpha_1$  and  $\alpha_2$  on  $\Delta$ , it follows that  $-\frac{1}{2} \log \sigma_a^2(\alpha)$  is bounded from above for any  $\alpha \in \Delta$ . Therefore, the statement that  $\sigma_a^2(\alpha)$  is uniformly bounded away from zero on  $\Delta$  can be established by a counter argument.

To show the uniqueness condition (A.1), write

$$\begin{aligned} n^{-1} [Q(\alpha) - Q(\alpha_0)] &= n^{-1} [Q_a(\alpha) - Q_a(\alpha_0)] + 2^{-1} [\log \sigma_a^2(\alpha) - \log \bar{\sigma}^2(\alpha)] \\ &\quad + 2^{-1} [\log \bar{\sigma}^2(\alpha_0) - \log \sigma_0^2], \end{aligned}$$

it follows, by Lemmas 4(1) and 6(1) and  $\bar{\sigma}^2(\alpha_0)$  being bounded away from zero, that  $\log \bar{\sigma}^2(\alpha_0) - \log \sigma_0^2 = o(1)$ . Moreover, we have already shown in (A.2) that  $\lim_{n \rightarrow \infty} [\sigma_a^2(\alpha) - \bar{\sigma}^2(\alpha)] \leq 0$ , hence,

$$\limsup_{n \rightarrow \infty} \max_{\alpha \in N^c(\alpha_0, \delta)} n^{-1} [Q(\alpha) - Q(\alpha_0)] \leq 0 \quad \text{for any } \delta > 0.$$

Now we will show that the above inequality holds strictly. It can be shown that  $n^{-1} Q(\alpha)$  is uniformly equicontinuous in  $\alpha$  on  $\Delta$  by Lemmas 4(1)–(3), 6 and 7(2) and the mean value theory. By the compactness of  $N^c(\alpha_0, \delta)$ , there exists an  $\delta > 0$  and a sequence  $\{\alpha_n\}$  in  $N^c(\alpha_0, \delta)$  converging to a point  $\alpha^* \neq \alpha_0$  such that  $\lim_{n \rightarrow \infty} n^{-1} [Q(\alpha_n) - Q(\alpha_0)] = 0$ . Because  $\lim_{n \rightarrow \infty} n^{-1} [Q(\alpha_n) - Q(\alpha^*)] = 0$  as  $\alpha_n \rightarrow \alpha^*$ , it follows that

$$(A.3) \quad \lim_{n \rightarrow \infty} n^{-1} [Q(\alpha^*) - Q(\alpha_0)] = 0.$$

Since  $Q_a(\alpha^*) - Q_a(\alpha_0) \leq 0$  and  $\lim_{n \rightarrow \infty} [\sigma_a^2(\alpha^*) - \bar{\sigma}^2(\alpha^*)] \leq 0$ , (A.3) is possible only if (i)  $\lim_{n \rightarrow \infty} [\sigma_a^2(\alpha^*) - \bar{\sigma}^2(\alpha^*)] = 0$  and (ii)  $\lim_{n \rightarrow \infty} n^{-1} [Q_a(\alpha^*) - Q_a(\alpha_0)] = 0$ . However, (i) is a contradiction when condition (7) holds by Lemmas 4(1)–(3), 6 and 7(2). If condition (7) holds, the contradiction follows from (ii) by condition (8).

The consistency of  $\hat{\sigma}^2$  can be obtained straightforwardly by Lemmas 4(1)–(3), 6, 7, 8(1)–(3), Chebyshev's inequality and  $\hat{\alpha} \xrightarrow{P} \alpha_0$ .  $\square$

PROOF OF THEOREM 2. Denoting  $\boldsymbol{\theta} = (\alpha, \sigma^2)^T$  and  $\boldsymbol{\theta}_0 = (\alpha_0, \sigma_0^2)^T$ , we get by Taylor's expansion that

$$0 = \frac{\partial l(\hat{\boldsymbol{\theta}})}{\partial \boldsymbol{\theta}} = \frac{\partial l(\boldsymbol{\theta}_0)}{\partial \boldsymbol{\theta}} + \frac{\partial^2 l(\tilde{\boldsymbol{\theta}})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^T} (\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0),$$

where  $\tilde{\boldsymbol{\theta}} = (\tilde{\alpha}, \tilde{\sigma}^2)^T$  lies between  $\hat{\boldsymbol{\theta}}$  and  $\boldsymbol{\theta}_0$ , and thus converges to  $\boldsymbol{\theta}_0$  in probability by Theorem 1. The asymptotic distribution of  $\hat{\boldsymbol{\theta}}$  can be obtained by showing that  $-\frac{1}{n} \frac{\partial^2 l(\tilde{\boldsymbol{\theta}})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^T} \xrightarrow{P} \Omega$  and  $\frac{1}{\sqrt{n}} \frac{\partial l(\boldsymbol{\theta}_0)}{\partial \boldsymbol{\theta}} \xrightarrow{D} N(\mathbf{0}, \Sigma + \Omega)$ , where  $\Omega$  is a nonsingular matrix by Theorem 1.

By straightforward calculations, it can be easily obtained that

$$\begin{aligned} \frac{1}{n} \frac{\partial^2 l(\boldsymbol{\theta})}{\partial \alpha^2} &= -\frac{1}{n} \text{tr}([WA^{-1}(\alpha)]^2) - \frac{1}{\sigma^2 n} (WY)^T P W Y, \\ \text{(A.4)} \quad \frac{1}{n} \frac{\partial^2 l(\boldsymbol{\theta})}{\partial \sigma^2 \partial \sigma^2} &= \frac{1}{2\sigma^4} - \frac{1}{\sigma^6 n} (A(\alpha)Y)^T P A(\alpha)Y, \\ \frac{1}{n} \frac{\partial^2 l(\boldsymbol{\theta})}{\partial \alpha \partial \sigma^2} &= -\frac{1}{\sigma^4 n} (WY)^T P A(\alpha)Y. \end{aligned}$$

As  $A(\tilde{\alpha})A^{-1} = I_n + (\alpha_0 - \tilde{\alpha})G$  by  $G = WA^{-1}$ , we have  $\frac{1}{n} \frac{\partial^2 l(\tilde{\boldsymbol{\theta}})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^T} - \frac{1}{n} \frac{\partial^2 l(\boldsymbol{\theta}_0)}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^T} = o_P(1)$  by Lemmas 6, 7, 8(1)–(3), Chebyshev's inequality, mean value theorem and  $\tilde{\boldsymbol{\theta}} \xrightarrow{P} \boldsymbol{\theta}_0$ . Furthermore, we have, by Lemmas 6, 7, 8(1)–(3) and Chebyshev's inequality that  $-\frac{1}{n} \frac{\partial^2 l(\boldsymbol{\theta}_0)}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^T} \xrightarrow{P} \Omega$ .

In the following, we will establish the asymptotic distribution of  $\frac{1}{\sqrt{n}} \frac{\partial l(\boldsymbol{\theta}_0)}{\partial \boldsymbol{\theta}}$ . It follows by Lemma 5(2) that  $\frac{1}{\sqrt{n}} (G\mathbf{m})^T P \mathbf{m} = o_P(1)$  when  $nh^6 \rightarrow 0$  and  $h^2 \log n \rightarrow 0$ . So, we have, by straightforward calculations, Lemmas 6(1) and 8, that

$$\frac{1}{\sqrt{n}} \frac{\partial l(\boldsymbol{\theta}_0)}{\partial \alpha} = \frac{1}{\sigma_0^2 \sqrt{n}} [(G\mathbf{m} - SG\mathbf{m})^T \boldsymbol{\varepsilon} + \{\boldsymbol{\varepsilon}^T G \boldsymbol{\varepsilon} - \sigma_0^2 \text{tr}(G)\}] + o_P(1)$$

and

$$\frac{1}{\sqrt{n}} \frac{\partial l(\boldsymbol{\theta}_0)}{\partial \sigma^2} = \frac{1}{2\sigma_0^4 \sqrt{n}} \{\boldsymbol{\varepsilon}^T \boldsymbol{\varepsilon} - n\sigma_0^2\} + o_P(1).$$



By straightforward calculations and Lemma 7(2), we have  $E(\frac{1}{n} \frac{\partial l(\theta_0)}{\partial \theta} \frac{\partial l(\theta_0)}{\partial \theta^T}) = \Sigma + \Omega + o(1)$ .

Finally, as the components of  $\frac{1}{\sqrt{n}} \frac{\partial l(\theta_0)}{\partial \theta} = (\frac{1}{\sqrt{n}} \frac{\partial l(\theta_0)}{\partial \alpha}, \frac{1}{\sqrt{n}} \frac{\partial l(\theta_0)}{\partial \sigma^2})^T$  are linear-quadratic forms of double arrays, using Lemma 9 we have  $\frac{1}{\sqrt{n}} \frac{\partial l(\theta_0)}{\partial \theta} \xrightarrow{D} N(\mathbf{0}, \Sigma + \Omega)$ .  $\square$

PROOF OF THEOREM 3. It can be easily shown that

$$\begin{aligned} & \sqrt{nh_1^2 f(s)} (\hat{\beta}(s) - \beta(s)) \\ &= \sqrt{nh_1^2 f(s)} (I_p, \mathbf{0}_{p \times 2p}) (\mathcal{X}_1^T \mathcal{W}_1 \mathcal{X}_1)^{-1} \mathcal{X}_1^T \mathcal{W}_1 \boldsymbol{\varepsilon} \\ & \quad + \sqrt{nh_1^2 f(s)} (\alpha_0 - \hat{\alpha}) (I_p, \mathbf{0}_{p \times 2p}) (\mathcal{X}_1^T \mathcal{W}_1 \mathcal{X}_1)^{-1} \mathcal{X}_1^T \mathcal{W}_1 WY \\ & \quad + \sqrt{nh_1^2 f(s)} \{ (I_p, \mathbf{0}_{p \times 2p}) (\mathcal{X}_1^T \mathcal{W}_1 \mathcal{X}_1)^{-1} \mathcal{X}_1^T \mathcal{W}_1 \mathbf{m} - \beta(s) \} \\ & \equiv J_{n1} + J_{n2} + J_{n3}, \end{aligned}$$

where  $\mathcal{X}_1$  and  $\mathcal{W}_1$  are  $\mathcal{X}$  and  $\mathcal{W}$  with  $h$  being replaced by  $h_1$ .

Let  $H_1$  be  $H$  with  $h$  being replaced by  $h_1$ . It follows by straightforward calculations that

$$\begin{aligned} & n^{-1} h_1^2 f(s) \text{cov}\{H_1^{-1} \mathcal{X}_1^T \mathcal{W}_1 \boldsymbol{\varepsilon}\} \\ &= \sigma_0^2 n^{-1} h_1^2 f(s) E\{H_1^{-1} \mathcal{X}_1^T \mathcal{W}_1^2 \mathcal{X}_1 H_1^{-1}\} \\ &= \sigma_0^2 f^2(s) \begin{pmatrix} \nu_0 \Psi + o_P(\mathbf{1}_p \mathbf{1}_p^T) & o_P(\mathbf{1}_p \mathbf{1}_{2p}^T) \\ o_P(\mathbf{1}_{2p} \mathbf{1}_p^T) & \nu_2 \Psi \otimes I_2 + o_P(\mathbf{1}_{2p} \mathbf{1}_{2p}^T) \end{pmatrix} \end{aligned}$$

this together with the central limit theorem, Lemma 2(1) and Slutsky's theorem lead to

$$J_{n1} \xrightarrow{D} N(\mathbf{0}, \nu_0 \kappa_0^{-2} \sigma_0^2 \Psi^{-1}).$$

It follows immediately from Lemmas 3, 2(1) and condition (4) that

$$(I_p, \mathbf{0}_{p \times 2p}) (\mathcal{X}_1^T \mathcal{W}_1 \mathcal{X}_1)^{-1} \mathcal{X}_1^T \mathcal{W}_1 G(\mathbf{m} + \boldsymbol{\varepsilon}) = o_P(1).$$

When  $nh_1^6 = O(1)$  and  $h/h_1 \rightarrow 0$ , we have  $\sqrt{\frac{h_1^2}{n}} (G\mathbf{m})^T P\mathbf{m} = o_P(1)$  using Lemma 5(1). It can be seen in the proof of Theorem 2 that  $\sqrt{nh_1^2} (\hat{\alpha} - \alpha_0) = o_P(1)$  under the assumptions of Theorem 3. Therefore,  $J_{n2} = o_P(1)$ .

The results of  $J_{n1}$  and  $J_{n2}$  together with Lemma 2(2),  $nh_1^6 = O(1)$  and  $h/h_1 \rightarrow 0$  lead to the theorem.  $\square$

PROOF OF THEOREM 4. It is obvious from the proof of nonsingularity of  $\Omega$  in Theorem 1 that  $\Omega$  is singular under condition (9).

Like Lee [9], to prove the consistency of  $\hat{\alpha}$ , it suffices to show that

$$\frac{\rho_n}{n} \{l_c(\alpha) - l_c(\alpha_0) - [Q(\alpha) - Q(\alpha_0)]\} = o_P(1) \quad \text{uniformly on } \Delta,$$

where  $Q(\alpha) = -n/2 \cdot \log \bar{\sigma}^2(\alpha) + \log |A(\alpha)|$  and  $\alpha_0$  is the unique maximizer.

It follows by the mean value theorem that

$$\begin{aligned} \frac{\rho_n}{n} \{l_c(\alpha) - l_c(\alpha_0) - [Q(\alpha) - Q(\alpha_0)]\} \\ = \frac{1}{\tilde{\sigma}^2(\tilde{\alpha})} \frac{\rho_n}{n} \left\{ [(WY)^T P A(\tilde{\alpha}) Y - L_n(\tilde{\alpha})] - \frac{\tilde{\sigma}^2(\tilde{\alpha}) - \bar{\sigma}^2(\tilde{\alpha})}{\tilde{\sigma}^2(\tilde{\alpha})} L_n(\tilde{\alpha}) \right\} \\ \times (\alpha - \alpha_0), \end{aligned}$$

where  $\tilde{\alpha}$  lies between  $\alpha$  and  $\alpha_0$ , and  $L_n(\tilde{\alpha}) = E[(WY)^T P A(\tilde{\alpha}) Y]$ . By the same arguments as in the proof of Theorem 1, we have  $\tilde{\sigma}^2(\tilde{\alpha}) - \bar{\sigma}^2(\tilde{\alpha}) = o_P(1)$  for any  $\tilde{\alpha}$  on  $\Delta$ , and  $\bar{\sigma}^2(\alpha)$  is uniformly bounded away from zero on  $\Delta$ . So,  $\tilde{\sigma}^2(\alpha)$  is uniformly bounded away from zero in probability. This together with Lemmas 4(5), 4(6), 11 and Chebyshev's inequality lead to

$$\frac{\rho_n}{n} \{l_c(\alpha) - l_c(\alpha_0) - [Q(\alpha) - Q(\alpha_0)]\} = o_P(1) \quad \text{uniformly on } \Delta.$$

The uniqueness condition of  $\alpha_0$  can be obtained by Lemma 4, Lemma 11, and the same arguments as in the proof of Theorem 1.  $\square$

**PROOF OF THEOREM 5.** By Taylor's expansion, we have that

$$0 = \frac{\partial l_c(\hat{\alpha})}{\partial \alpha} = \frac{\partial l_c(\alpha_0)}{\partial \alpha} + \frac{\partial^2 l_c(\tilde{\alpha})}{\partial \alpha^2} (\hat{\alpha} - \alpha_0),$$

where  $\tilde{\alpha}$  lies between  $\hat{\alpha}$  and  $\alpha_0$ , and thus converges to  $\alpha_0$  in probability by Theorem 4. So, the asymptotic distribution of  $\hat{\alpha}$  can be obtained by proving that

$$-\frac{\rho_n}{n} \frac{\partial^2 l_c(\tilde{\alpha})}{\partial \alpha^2} \xrightarrow{P} \sigma_1^2 \quad \text{and} \quad \sqrt{\frac{\rho_n}{n}} \frac{\partial l_c(\alpha_0)}{\partial \alpha} \xrightarrow{D} N(0, \sigma_2^2 / \sigma_0^4),$$

when  $\rho_n \rightarrow \infty$ , where  $\sigma_1^2 = \frac{1}{\sigma_0^2} \lim_{n \rightarrow \infty} \frac{\rho_n}{n} E[(G\mathbf{m} - SG\mathbf{m})^T (G\mathbf{m} - SG\mathbf{m})]$  and  $\sigma_2^2 = \sigma_0^4 \sigma_1^2$ .

As we have, by  $A(\alpha)A^{-1} = I_n + (\alpha_0 - \alpha)G$ , Lemma 11 and Chebyshev's inequality, that  $\frac{\rho_n}{n} (WY)^T P WY = O_P(1)$  and  $\frac{\rho_n}{n} (WY)^T P A(\alpha)Y = O_P(1)$ , so, when  $\rho_n \rightarrow \infty$ ,

$$\frac{\rho_n}{n} \frac{\partial^2 l_c(\alpha)}{\partial \alpha^2} = -\frac{1}{\tilde{\sigma}^2(\alpha)} \cdot \frac{\rho_n}{n} (WY)^T P WY - \frac{\rho_n}{n} \text{tr}([WA^{-1}(\alpha)]^2) + o_P(1).$$

This together with Lemmas 6(1), 8(1) lead to  $\tilde{\sigma}^2(\alpha) = \sigma_0^2 + o_P(1)$  for any  $\alpha \in \Delta$  when  $\rho_n \rightarrow \infty$ . Therefore, by the mean value theorem, conditions (5)–(6) and  $\tilde{\alpha} \xrightarrow{P} \alpha_0$ , we have  $\frac{\rho_n}{n} \left\{ \frac{\partial^2 l_c(\tilde{\alpha})}{\partial \alpha^2} - \frac{\partial^2 l_c(\alpha_0)}{\partial \alpha^2} \right\} = o_P(1)$ .

It follows, from  $\tilde{\sigma}^2(\alpha_0) \xrightarrow{P} \sigma_0^2$ , Lemma 11, Chebyshev's inequality and the row sums of  $G$  being uniform order  $O(1/\sqrt{\rho_n})$ , that  $-\frac{\rho_n}{n} \frac{\partial^2 l_c(\alpha_0)}{\partial \alpha^2} \xrightarrow{P} \sigma_1^2$ .

In the following, we will establish the asymptotic distribution of  $\sqrt{\frac{\rho_n}{n}} \frac{\partial l_c(\alpha_0)}{\partial \alpha}$ .

By Lemmas 10(2) and 11(3), it is easy to see  $\sqrt{\frac{\rho_n}{n}} (G\mathbf{m})^T P\mathbf{m} = o_P(1)$  and  $\sqrt{\frac{\rho_n}{n}} (G\boldsymbol{\varepsilon})^T P\mathbf{m} = o_P(1)$  when  $nh^6 \rightarrow 0$  and  $h^2 \log n \rightarrow 0$ . By straightforward calculations and Lemmas 6(1), 8(1), 11(5) and 11(7), we have the first-order derivative of  $\sqrt{\frac{\rho_n}{n}} l_c(\alpha)$  at  $\alpha_0$  is

$$\frac{1}{\tilde{\sigma}^2(\alpha_0)} \sqrt{\frac{\rho_n}{n}} \left\{ (G\mathbf{m} - SG\mathbf{m})^T \boldsymbol{\varepsilon} + \boldsymbol{\varepsilon}^T \left[ G - \frac{1}{n} \text{tr}(G) I_n \right] \boldsymbol{\varepsilon} \right\} + o_P(1).$$

By Lemma 12, we have

$$\sigma_{qn}^{-1} \left\{ (G\mathbf{m} - SG\mathbf{m})^T \boldsymbol{\varepsilon} + \boldsymbol{\varepsilon}^T \left[ G - \frac{1}{n} \text{tr}(G) I_n \right] \boldsymbol{\varepsilon} \right\} \xrightarrow{D} N(0, 1),$$

where  $\sigma_{qn}^2 = \text{var}\{(G\mathbf{m} - SG\mathbf{m})^T \boldsymbol{\varepsilon} + \boldsymbol{\varepsilon}^T [G - \frac{1}{n} \text{tr}(G) I_n] \boldsymbol{\varepsilon}\}$ . So, by  $\frac{\rho_n}{n} \sigma_{qn}^2 \rightarrow \sigma_2^2$  and  $\tilde{\sigma}^2(\alpha_0) \xrightarrow{P} \sigma_0^2$ , we have  $\sqrt{\frac{n}{\rho_n}} (\hat{\alpha} - \alpha_0) \xrightarrow{D} N(0, \sigma_0^2 \lambda_4^{-1})$ .  $\square$

**PROOF OF THEOREM 6.** By straightforward calculations, Lemmas 6(1), 8(1), 11, Chebyshev's inequality and Theorem 5, we have  $\sqrt{n}(\hat{\sigma}^2 - \sigma_0^2) = \frac{1}{\sqrt{n}} \sum_{i=1}^n (\varepsilon_i^2 - \sigma_0^2) + o_P(1)$  when  $\rho_n \rightarrow \infty$ . This together with the central limit theorem lead to Theorem 6.  $\square$

**PROOF OF THEOREM 7.** Theorem 7 can be obtained by using the same arguments as in the proof of Theorem 3, except that here

$$\begin{aligned} J_{n2} &= \sqrt{f(s)} \sqrt{\frac{nh_1^2}{\rho_n}} (\alpha_0 - \hat{\alpha}) (I_p, \mathbf{0}_{p \times 2p}) (n^{-1} H_1^{-1} \mathcal{X}_1^T \mathcal{W}_1 \mathcal{X}_1 H_1^{-1})^{-1} \\ &\quad \times \frac{\sqrt{\rho_n}}{n} H_1^{-1} \mathcal{X}_1^T \mathcal{W}_1 G(\mathbf{m} + \boldsymbol{\varepsilon}). \end{aligned}$$

By Lemma 2(1), Markov's inequality, the row sums of the matrix  $G$  having uniform order  $O(1/\sqrt{\rho_n})$  and condition (4), we have

$$(I_p, \mathbf{0}_{p \times 2p}) (n^{-1} H_1^{-1} \mathcal{X}_1^T \mathcal{W}_1 \mathcal{X}_1 H_1^{-1})^{-1} \frac{\sqrt{\rho_n}}{n} H_1^{-1} \mathcal{X}_1^T \mathcal{W}_1 G(\mathbf{m} + \boldsymbol{\varepsilon}) = O_P(1).$$

Furthermore, it can be seen from the proof of Theorem 5 and Lemma 10(1) that when  $nh_1^6 = O(1)$  and  $h/h_1 \rightarrow 0$ ,  $\sqrt{\frac{nh_1^2}{\rho_n}} (\hat{\alpha} - \alpha) \xrightarrow{P} 0$ . So,  $J_{n2} = o_P(1)$ .  $\square$

## SUPPLEMENTARY MATERIAL

**Detailed proofs of lemmas and theorems** (DOI: [10.1214/13-AOS1201SUPP](https://doi.org/10.1214/13-AOS1201SUPP); .pdf). We provide the detailed proofs of the lemmas and theorems.

## REFERENCES

- [1] ANSELIN, L. (1988). *Spatial Econometrics: Methods and Models*. Kluwer Academic, Dordrecht.
- [2] CHENG, M.-Y., ZHANG, W. and CHEN, L.-H. (2009). Statistical estimation in generalized multiparameter likelihood models. *J. Amer. Statist. Assoc.* **104** 1179–1191. [MR2750243](#)
- [3] FAN, J. and GIJBELS, I. (1996). *Local Polynomial Modelling and Its Applications*. Chapman & Hall, London. [MR1383587](#)
- [4] FAN, J. and ZHANG, W. (1999). Statistical estimation in varying coefficient models. *Ann. Statist.* **27** 1491–1518. [MR1742497](#)
- [5] FAN, J. and ZHANG, W. (2000). Simultaneous confidence bands and hypothesis testing in varying-coefficient models. *Scand. J. Stat.* **27** 715–731. [MR1804172](#)
- [6] GAO, J., LU, Z. and TJØSTHEIM, D. (2006). Estimation in semiparametric spatial regression. *Ann. Statist.* **34** 1395–1435. [MR2278362](#)
- [7] HANSEN, B. E. (2008). Uniform convergence rates for kernel estimation with dependent data. *Econometric Theory* **24** 726–748. [MR2409261](#)
- [8] KELEJIAN, H. H. and PRUCHA, I. R. (2010). Specification and estimation of spatial autoregressive models with autoregressive and heteroskedastic disturbances. *J. Econometrics* **157** 53–67. [MR2652278](#)
- [9] LEE, L.-F. (2004). Asymptotic distributions of quasi-maximum likelihood estimators for spatial autoregressive models. *Econometrica* **72** 1899–1925. [MR2095537](#)
- [10] LI, J. and ZHANG, W. (2011). A semiparametric threshold model for censored longitudinal data analysis. *J. Amer. Statist. Assoc.* **106** 685–696. [MR2847950](#)
- [11] LINTON, O. (1995). Second order approximation in the partially linear regression model. *Econometrica* **63** 1079–1112. [MR1348514](#)
- [12] ORD, K. (1975). Estimation methods for models of spatial interaction. *J. Amer. Statist. Assoc.* **70** 120–126. [MR0413393](#)
- [13] SU, L. and JIN, S. (2010). Profile quasi-maximum likelihood estimation of partially linear spatial autoregressive models. *J. Econometrics* **157** 18–33. [MR2652276](#)
- [14] SUN, Y., YAN, H., ZHANG, W. and LU, Z. (2014). Supplement to “A semiparametric spatial dynamic model.” DOI:[10.1214/13-AOS1201SUPP](https://doi.org/10.1214/13-AOS1201SUPP).
- [15] SUN, Y., ZHANG, W. and TONG, H. (2007). Estimation of the covariance matrix of random effects in longitudinal studies. *Ann. Statist.* **35** 2795–2814. [MR2382666](#)
- [16] TAO, H. and XIA, Y. (2012). Adaptive semi-varying coefficient model selection. *Statist. Sinica* **22** 575–599. [MR2954353](#)
- [17] WANG, H. and XIA, Y. (2009). Shrinkage estimation of the varying coefficient model. *J. Amer. Statist. Assoc.* **104** 747–757. [MR2541592](#)
- [18] WHITE, H. (1994). *Estimation, Inference and Specification Analysis*. Cambridge Univ. Press, Cambridge. [MR1292251](#)
- [19] ZHANG, W., FAN, J. and SUN, Y. (2009). A semiparametric model for cluster data. *Ann. Statist.* **37** 2377–2408. [MR2543696](#)

- [20] ZHANG, W., LEE, S.-Y. and SONG, X. (2002). Local polynomial fitting in semivarying coefficient model. *J. Multivariate Anal.* **82** 166–188. [MR1918619](#)

Y. SUN  
SCHOOL OF ECONOMICS  
SHANGHAI UNIVERSITY OF FINANCE AND ECONOMICS  
NO. 777 GUODING ROAD  
SHANGHAI, 200433  
P.R. CHINA  
E-MAIL: [sunyan@mail.shufe.edu.cn](mailto:sunyan@mail.shufe.edu.cn)

H. YAN  
W. ZHANG  
DEPARTMENT OF MATHEMATICS  
THE UNIVERSITY OF YORK  
HESLINGTON  
YORK, YO10 5DD  
UNITED KINGDOM  
E-MAIL: [yanli8626@qq.com](mailto:yanli8626@qq.com)  
[wenyang.zhang@york.ac.uk](mailto:wenyang.zhang@york.ac.uk)

Z. LU  
SCHOOL OF MATHEMATICAL SCIENCES  
THE UNIVERSITY OF SOUTHAMPTON  
HIGHFIELD  
SOUTHAMPTON, SO17 1BJ  
UNITED KINGDOM  
E-MAIL: [Z.Lu@soton.ac.uk](mailto:Z.Lu@soton.ac.uk)