

Semiparametric method and theory for continuously indexed spatio-temporal processes

Jialuo Liu^a, Tingjin Chu^{b,*}, Jun Zhu^c, Haonan Wang^a

^a Department of Statistics, Colorado State University, United States of America

^b School of Mathematics and Statistics, University of Melbourne, Australia

^c Department of Statistics, University of Wisconsin-Madison, United States of America

ARTICLE INFO

Article history:

Received 19 February 2020

Received in revised form 4 February 2021

Accepted 5 February 2021

Available online 13 February 2021

AMS 2010 subject classifications:

primary 62H11

secondary 62M30

Keywords:

Bimodal kernel

Random fields

Spatial statistics

Spatio-temporal statistics

ABSTRACT

Spatio-temporal processes with a continuous index in space and time are useful for modeling spatio-temporal data in many scientific disciplines such as environmental and health sciences. However, approaches that enable simultaneous estimation of the mean and covariance functions for such spatio-temporal processes are limited. Here, we propose a flexible spatio-temporal model with partially linear regression in the mean function and local stationarity in the covariance function. We develop a profile likelihood method for estimation and an effective bandwidth selection in the presence of spatio-temporally correlated errors. Specifically, we employ a family of bimodal kernels to alleviate bias, which may be of independent interest for semiparametric spatial statistics. The theoretical properties of our profile likelihood estimation, including consistency and asymptotic normality, are established. A simulation study is conducted and suggests sound empirical properties, while a health hazard data example further illustrates the methodology.

© 2021 Elsevier Inc. All rights reserved.

1. Introduction

In this paper, we develop new semiparametric methodology and theory for spatio-temporal processes where both space and time are continuously indexed, which often arise in many scientific disciplines [see, e.g., 34]. An illustrative data set comprises measurements of a health hazard taken in an indoor environment by both static sensors at fixed sampling locations and roving sensors at varying sampling locations over time [31]. The spatio-temporal sampling design is non-standard due to data irregularity and sparsity in both space and time, calling for development of novel methodology and theory.

In spatial statistics, geostatistical data with continuous spatial index and lattice data with discrete spatial index often require different modeling techniques. For example, to account for spatial dependence, a Matérn covariance function is typically used for geostatistical data, while a spatial weight matrix is used for lattice data [8]. For spatio-temporal datasets, the distinction between continuous and discrete index applies to both spatial and temporal dimensions. To analyze datasets with continuous spatial index and discrete temporal index, time series methods for temporal data are often combined with geostatistical methods for spatial data. For example, Stroud et al. [41] developed a state space model where spatial variability is captured by a locally weighted mixture of linear regressions while the regression coefficients are allowed to vary with time. Sun et al. [42] proposed a profile likelihood based estimation procedure for a semiparametric

* Corresponding author.

E-mail address: tingjin.chu@unimelb.edu.au (T. Chu).

spatial dynamic model with a nonlinear spatial trend. Al-Sulami et al. [1] considered a nonlinear spatio-temporal model to investigate the relationship between housing price index (HPI) and consumer price index (CPI) for individual states in the USA. The aforementioned spatial time series methods can capture nonlinearity and nonstationarity in space and/or time, assuming that data are observed at regular and discrete time points. In contrast, for spatio-temporal data with continuous temporal index, approaches that enable simultaneous estimation of the mean and covariance functions are limited. While the existing methods focus primarily on linear regression models [see, e.g., 11], we will develop semiparametric methods and theory for continuously indexed spatio-temporal processes.

The underlying spatio-temporal process can be decomposed into a mean trend and a spatio-temporal error process. Partially linear models offer a flexible way to model the mean trend of spatio-temporal data. For independent data, partially linear models have been extensively studied [see, e.g., 14,22,26,27,37]. For spatio-temporal data, Gao et al. [18] proposed an estimation procedure based on marginal integration for geostatistical partially linear models, and Lu et al. [29] developed spatio-temporal varying coefficient models, which can be applied to spatio-temporal partially linear models. Theoretical property is established for both works under the spatio-temporal mixing conditions. Since both works focus on estimating the mean trend of spatio-temporal data, spatio-temporal error is treated as independent in estimation. In practice, there is considerable interest in spatio-temporal covariance functions, which characterize the spatio-temporal dependence of underlying processes. Furthermore, to interpolate unsampled spatial locations and time points (spatio-temporal kriging), spatio-temporal covariance functions are a key building block. Thus, there is clearly a need for statistical methods to estimate spatio-temporal covariance functions and here, we aim to develop new methodology which allows simultaneous estimation of the mean and covariance functions.

Various types of spatio-temporal covariance functions have been developed [see, e.g., 7,9,16,19,35,39]. However, the dependence structure in spatio-temporal data poses challenges for establishing the asymptotic properties. In spatial statistics, there are three commonly used asymptotic frameworks, namely, increasing-domain asymptotics, fixed-domain asymptotics and mixed-domain asymptotics. For increasing-domain asymptotics, the spatial domain expands as the number of observations increases [see, e.g., 6,10,32]. For fixed-domain asymptotics, the spatial domain is fixed and the sampling locations get denser [see, e.g., 28,38,45,46]. A mixed-domain asymptotic framework allows both spatial domain and sampling density to increase [see, e.g., 4,21,23,30]. For spatio-temporal processes, Bandyopadhyay et al. [3] considered an increasing temporal domain and a mixed spatial domain for a Fourier analysis. Chu et al. [5] proposed a spatio-temporal expanding distance (STED) asymptotic framework in a fixed spatio-temporal domain, which extends the aforementioned asymptotic frameworks for spatial domain to spatio-temporal domain for exploring the asymptotic properties of statistical inference for spatio-temporal processes. The STED framework also paves the way for studying the local behavior of a spatio-temporal process, especially the second-order properties. Here, we will consider a locally stationary spatio-temporal covariance function, introduced by Chu et al. [5], to study the slowly-varying second-order nonstationarity under the STED asymptotic framework.

In essence, the mean trend of spatio-temporal data is modeled through partially linear models, and spatio-temporal dependence is accounted by locally stationary spatio-temporal covariance functions. The resulting model provides a flexible way to analyze continuously indexed spatio-temporal datasets. For estimation, the main challenge is to incorporate spatio-temporal covariance functions, which we overcome by profiling the spatio-temporal likelihood function. In addition, the theoretical property of proposed method is investigated under STED framework, and both consistency and asymptotic normality are established. Furthermore, a proper bandwidth selection is critical for estimation. For *iid* data, various methods have been studied, notably cross validation [see, e.g., 13], but are known to not perform well for non-*iid* data [see, e.g., 2,12,24,33]. Here, we show that cross-validation is asymptotically biased in the presence of spatio-temporal correlated errors for most commonly used kernels. We also propose a cross-validation based method with bimodal kernels to alleviate this bias in bandwidth selection.

The remainder of the paper is organized as follows. Section 2 introduces the spatio-temporal model and the profile likelihood method. The asymptotic properties of the profile likelihood estimation are established in Section 3 under suitable regularity conditions. In Section 4, we discuss the choice of kernel functions and develop a procedure for bandwidth selection. Numerical examples including a simulation study and the health hazard data example are given in Sections 5 and 6, respectively. Appendix A contains the technical details including proofs, while additional simulation results are given as Supplementary Materials.

2. Model and estimation

2.1. Spatio-temporal semiparametric model

We consider the following spatio-temporal process for the response variable $y(\cdot, \cdot)$,

$$y(\mathbf{s}, t) = \mathbf{x}(\mathbf{s}, t)^\top \boldsymbol{\beta} + f(t) + \varepsilon(\mathbf{s}, t), \quad \mathbf{s} \in \mathcal{R}, \quad t \in \mathcal{T}, \quad (1)$$

where the location \mathbf{s} is in the unit hypercube $\mathcal{R} = [0, 1]^d$ for $d \geq 1$ and the rescaled time t takes values in $\mathcal{T} = [0, 1]$. Here, $\mathbf{x}(\mathbf{s}, t) = (x_1(\mathbf{s}, t), \dots, x_p(\mathbf{s}, t))^\top$ is a $p \times 1$ vector of covariates at spatial location \mathbf{s} and time point t , $\boldsymbol{\beta} = (\beta_1, \dots, \beta_p)^\top$ is a $p \times 1$ vector of regression coefficients, and $f(t)$ denotes a fixed nonparametric temporal function. In the special case of $\boldsymbol{\beta} = \mathbf{0}$, (1) has a fully nonparametric mean function. Furthermore, the zero-mean spatio-temporal Gaussian random

process $\varepsilon(\mathbf{s}, t)$ accounts for the local variations unexplained by the mean function (i.e., trend) $\mathbf{x}(\mathbf{s}, t)^\top \boldsymbol{\beta} + f(t)$. In practice, $\varepsilon(\mathbf{s}, t) = \varepsilon_1(\mathbf{s}, t) + \varepsilon_2(\mathbf{s}, t)$, where $\varepsilon_1(\mathbf{s}, t)$ is a Gaussian spatio-temporal error process, and $\varepsilon_2(\mathbf{s}, t)$ are Gaussian iid errors with mean 0 and equal variance, representing the nugget effect and independent of $\varepsilon_1(\mathbf{s}, t)$. Denote $\gamma((\mathbf{s}, t), (\mathbf{s}', t'); \boldsymbol{\theta})$ as the covariance function of $\varepsilon(\mathbf{s}, t)$, where $(\mathbf{s}, t), (\mathbf{s}', t') \in \mathcal{R} \times \mathcal{T}$ and $\boldsymbol{\theta}$ is a $q \times 1$ vector of covariance function parameters. For $\gamma((\mathbf{s}, t), (\mathbf{s}', t'); \boldsymbol{\theta})$, various types of covariance functions are proposed [see, e.g., 7,9,16,19,35,39]. In this work, we focus on locally stationary covariance function, introduced by Chu et al. [5], and more details are provided in Section 5.

We consider N samples observed at $(\mathbf{s}_1, t_1), \dots, (\mathbf{s}_N, t_N)$. Define the $N \times 1$ vector of the response variable as $\mathbf{y} = (y(\mathbf{s}_1, t_1), \dots, y(\mathbf{s}_N, t_N))^\top$, the $N \times p$ design matrix as $\mathbf{X} = [x_j(\mathbf{s}_i, t_i)]_{i=1, j=1}^{N, p}$, and the $N \times 1$ vector of the errors as $\boldsymbol{\varepsilon} = (\varepsilon(\mathbf{s}_1, t_1), \dots, \varepsilon(\mathbf{s}_N, t_N))^\top$. Let $\mathbf{f} = (f(t_1), \dots, f(t_N))^\top$ denote the temporal function at the N sampling points. We have

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{f} + \boldsymbol{\varepsilon}. \quad (2)$$

The $N \times N$ covariance matrix of $\boldsymbol{\varepsilon}$ is expressed as $\boldsymbol{\Gamma}(\boldsymbol{\theta}) = [\gamma((\mathbf{s}_i, t_i), (\mathbf{s}_j, t_j); \boldsymbol{\theta})]_{i,j=1}^N$. Further, let $\boldsymbol{\eta} = (\boldsymbol{\beta}^\top, \boldsymbol{\theta}^\top)^\top$ denote a $(p+q) \times 1$ vector of parameters comprising both the regression coefficients $\boldsymbol{\beta}$ and the covariance function parameters $\boldsymbol{\theta}$.

2.2. Profile likelihood estimation

Since the likelihood principle cannot be easily adopted for semiparametric models like (1), here we develop a profile likelihood method for model estimation. For a given $\boldsymbol{\beta}$, let $y_i^* = y(\mathbf{s}_i, t_i) - \mathbf{x}(\mathbf{s}_i, t_i)^\top \boldsymbol{\beta}$ denote a partially detrended spatio-temporal process for the response variable and let $\mathbf{y}^* = (y_1^*, \dots, y_N^*)^\top$ denote an $N \times 1$ vector of partially detrended spatio-temporal response variables. We obtain an estimate of \mathbf{f} by local polynomial regression; that is, we minimize the following criterion, with respect to $\mathbf{b}_t = (b_{0,t}, b_{1,t})^\top$,

$$\sum_{i=1}^N \{y_i^* - b_{0,t} - b_{1,t}(t_i - t)\}^2 K_h(t_i - t), \quad (3)$$

where $K_h = K(\cdot/h)/h$ is a kernel function $K(\cdot)$ with a bandwidth h .

With $\mathbf{K}_t = \text{diag}\{K_h(t_1 - t), \dots, K_h(t_N - t)\}$, $\mathbf{D}_t = (\mathbf{1}_N, \mathbf{d}_{1t})$, $\mathbf{d}_{1t} = ((t_1 - t)/h, \dots, (t_N - t)/h)^\top$, and $\mathbf{1}_N$ is an $N \times 1$ vector of 1's, it follows from (1) that, $(\hat{b}_{0,t}, h\hat{b}_{1,t})^\top = \boldsymbol{\omega}(t)\mathbf{y}^*$, where $\boldsymbol{\omega}(t) = (\mathbf{D}_t^\top \mathbf{K}_t \mathbf{D}_t)^{-1} \mathbf{D}_t^\top \mathbf{K}_t$. The resulting estimate of \mathbf{f} is $\tilde{\mathbf{f}} = \mathbf{S}\mathbf{y}^* = \mathbf{S}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})$, where the smoother matrix is

$$\mathbf{S} = (\boldsymbol{\omega}_1(t_1)^\top, \dots, \boldsymbol{\omega}_1(t_N)^\top)^\top, \quad (4)$$

and $\boldsymbol{\omega}_1(t) = (1, 0)\boldsymbol{\omega}(t)$. Plugging $\tilde{\mathbf{f}}$ into (2), we have the following approximation

$$(\mathbf{I} - \mathbf{S})\mathbf{y} \approx (\mathbf{I} - \mathbf{S})\mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}. \quad (5)$$

If $\boldsymbol{\varepsilon}$ is a sequence of independent and identically distributed random variables, the profile method is used to obtain estimates of $\boldsymbol{\beta}$ as

$$\bar{\boldsymbol{\beta}} = \{\mathbf{X}^\top (\mathbf{I} - \mathbf{S})^\top (\mathbf{I} - \mathbf{S})\mathbf{X}\}^{-1} \mathbf{X}^\top (\mathbf{I} - \mathbf{S})^\top (\mathbf{I} - \mathbf{S})\mathbf{y}$$

and $\bar{\mathbf{f}} = \mathbf{S}(\mathbf{y} - \mathbf{X}\bar{\boldsymbol{\beta}})$. However, the above method does not account for the spatio-temporal dependence of $\varepsilon(\mathbf{s}, t)$, and therefore, the spatio-temporal parameter $\boldsymbol{\theta}$ cannot be estimated. In order to estimate both the mean trend and the spatio-temporal parameter $\boldsymbol{\theta}$, we propose to maximize the approximated profile log-likelihood function based on (5),

$$\ell(\boldsymbol{\beta}, \boldsymbol{\theta}) = -(N/2) \log(2\pi) - (1/2) \log\{\det \boldsymbol{\Gamma}(\boldsymbol{\theta})\} - (1/2) (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top (\mathbf{I} - \mathbf{S})^\top \boldsymbol{\Gamma}(\boldsymbol{\theta})^{-1} (\mathbf{I} - \mathbf{S})(\mathbf{y} - \mathbf{X}\boldsymbol{\beta}). \quad (6)$$

The estimate of $(\boldsymbol{\beta}^\top, \boldsymbol{\theta}^\top)^\top$ is the maximizer of (6), and is denoted as $(\hat{\boldsymbol{\beta}}^\top, \hat{\boldsymbol{\theta}}^\top)^\top$. Consequently, the estimate of \mathbf{f} can be expressed as

$$\hat{\mathbf{f}} = \mathbf{S}(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}).$$

In addition, let $\mathbf{f}' = (f'(t_1), \dots, f'(t_N))^\top$ denote an $N \times 1$ vector of the first-order derivatives of the temporal function $f(t)$ evaluated at the sampling time points $\{t_1, \dots, t_N\}$. Minimizing (3) yields an estimate of \mathbf{f}'

$$\hat{\mathbf{f}}' = \mathbf{L}(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}),$$

where $\mathbf{L} = (h^{-1}\boldsymbol{\omega}_2(t_1)^\top, \dots, h^{-1}\boldsymbol{\omega}_2(t_N)^\top)^\top$ and $\boldsymbol{\omega}_2(t) = (0, 1)\boldsymbol{\omega}(t)$.

In general, we write the estimate of $\mathbf{F}(t) = (f(t), hf'(t))^\top$ as $\hat{\mathbf{F}}(t) = \boldsymbol{\omega}(t)(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})$. In the case of spatio-temporal independence (i.e., $\boldsymbol{\Gamma} = \sigma^2 \mathbf{I}$), the estimates of $\boldsymbol{\beta}$ and σ^2 in (6) can be expressed in closed form [see, e.g., 14]. In the case of a nonparametric mean function (i.e., $\boldsymbol{\beta} = \mathbf{0}$), (6) can still be maximized to obtain the estimates of $\boldsymbol{\theta}$ and \mathbf{f} , while the estimate of \mathbf{f}' can be obtained by $\hat{\mathbf{f}}' = \mathbf{L}\mathbf{y}$. The estimation procedure above depends on the choice of bandwidth, which will be discussed in Section 4.

3. Asymptotic properties

3.1. Asymptotic framework

The spatio-temporal framework is essential for establishing asymptotic property of parameter estimators. In spatial statistics, there are three asymptotic frameworks, namely, increasing-domain asymptotics, fixed-domain asymptotics, and mixed-domain asymptotics. For the likelihood-based method considered here, we employ a spatio-temporal expanding distance (STED) asymptotic framework in a fixed spatio-temporal domain, which provides a flexible tool for exploring the asymptotic properties of statistical inference for spatio-temporal processes [5]. Let n denote the stage of the asymptotics and let $\{A_n\}$ and $\{B_n\}$ be two sequences of increasing positive numbers. The (A_n, B_n) -rate STED asymptotic framework in a fixed spatio-temporal domain is defined as follows. For all n , there exist positive constants c_1 , c_2 and c_3 , independent of n , such that

$$(A.1) \quad \delta_n / \min_{1 \leq j \leq N_n} \delta_{j,n} \leq c_1,$$

$$(A.2) \quad \zeta_n / \min_{1 \leq j \leq N_n} \zeta_{j,n} \leq c_2,$$

$$(A.3) \quad \delta_n^d A_n^d \zeta_n B_n \geq c_3,$$

where $\delta_{j,n} = \min\{\|\mathbf{s}_i - \mathbf{s}_j\| : 1 \leq i \leq N_n, \mathbf{s}_i \neq \mathbf{s}_j\}$, $\delta_n = \max_{1 \leq j \leq N_n} \delta_{j,n}$, $\zeta_{j,n} = \min\{|t_i - t_j| : 1 \leq i \leq N_n, t_i \neq t_j\}$ and $\zeta_n = \max_{1 \leq j \leq N_n} \zeta_{j,n}$. We assume that the error process $\varepsilon(\mathbf{s}, t)$ is locally stationary in the sense that a covariance function $\gamma_n((\mathbf{s}, t), (\mathbf{s}', t'))$ is said to be locally stationary if there exists a sequence of functions $g_n(\cdot, \cdot, \mathbf{s}, t)$ such that,

$$|\gamma_n((\mathbf{s}, t), (\mathbf{s}', t')) - g_n(\mathbf{s}' - \mathbf{s}, t' - t, \mathbf{s}, t)| = \mathcal{O}(\|\mathbf{s}' - \mathbf{s}\| + |t' - t| + \rho_n),$$

uniformly for all $(\mathbf{s}, t), (\mathbf{s}', t') \in \mathcal{R} \times \mathcal{T}$, where $\{\rho_n\}$ is a sequence of positive numbers, which does not depend on the location, time or the parameter $\boldsymbol{\theta}$. Furthermore, $\rho_n \rightarrow 0$ as $n \rightarrow \infty$. In addition, there exists a function g such that, as $n \rightarrow \infty$,

$$\lim_{n \rightarrow \infty} |g_n(\mathbf{s}' - \mathbf{s}, t' - t, \mathbf{s}, t) - g(\mathbf{u}_1, u_2, \mathbf{s}, t)| \rightarrow 0$$

uniformly for all $(\mathbf{s}, t), (\mathbf{s}', t') \in \mathcal{R} \times \mathcal{T}$, where $\mathbf{u}_1 = A_n(\mathbf{s}' - \mathbf{s})$ and $u_2 = B_n(t' - t)$.

We use a one-dimensional (1D) toy example to illustrate the structure of the locally stationary covariance function. For locations $s \in \mathcal{R} = [0, 1]$, we construct a locally stationary covariance function by taking the product of a positive function $D(s)$ and a stationary covariance function such that $\gamma(s, s') = D(s)D(s')\exp(-d/r)$, where r is the range parameter and $d = |s - s'|$ is the distance between s and s' . Fig. A of the supplementary material demonstrates four covariance functions, one stationary covariance function where $D_1(s) = 1$ and three locally stationary covariance functions.

3.2. Asymptotic properties

For iid data, the maximum profile likelihood estimate $\hat{\boldsymbol{\beta}}$ is consistent and asymptotically normal [14]. For the spatio-temporal semiparametric model (1), the asymptotic properties of the maximum profile likelihood estimates ${}^n\hat{\boldsymbol{\beta}}$ and ${}^n\hat{\boldsymbol{\theta}}$, which maximize (6), will be established as follows.

Theorem 1. Under (C.1)–(C.13) in the Appendix A, there exists, with probability tending to one, a local maximizer ${}^n\hat{\boldsymbol{\eta}} = ({}^n\hat{\boldsymbol{\beta}}^\top, {}^n\hat{\boldsymbol{\theta}}^\top)^\top$ of $\ell(\boldsymbol{\eta})$ such that $\|{}^n\hat{\boldsymbol{\eta}} - \boldsymbol{\eta}_0\| = \mathcal{O}_p(N_n^{-1/2})$. Moreover, the local maximizer ${}^n\hat{\boldsymbol{\eta}}$ is asymptotic normal; that is, as $n \rightarrow \infty$,

$$N_n^{1/2}({}^n\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}_0) \xrightarrow{D} \mathcal{N}(\mathbf{0}, \boldsymbol{\Pi}^{-1}) \quad \text{and} \quad N_n^{1/2}({}^n\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0) \xrightarrow{D} \mathcal{N}(\mathbf{0}, \boldsymbol{\mathcal{I}}_0(\boldsymbol{\theta}_0)^{-1}).$$

Theorem 1 establishes that the estimate ${}^n\hat{\boldsymbol{\eta}}$ is root- N_n consistent. However, the asymptotic variance of ${}^n\hat{\boldsymbol{\beta}}$ does not converge to the information matrix (13). As will be seen in Appendix A, if $\|\boldsymbol{\Gamma}^{-1}\|_\infty = \mathcal{O}(1)$, then $\mathbf{X}^\top \boldsymbol{\Gamma}^{-1} \mathbf{X} \succeq \boldsymbol{\Phi}^\top \boldsymbol{\Gamma}^{-1} \boldsymbol{\Phi}$, where $\mathbf{A} \succeq \mathbf{B}$ if $\mathbf{A} - \mathbf{B}$ is positive semi-definite. That is, the asymptotic variance of $\hat{\boldsymbol{\beta}}$ in partially linear models is greater than those in simple linear regression models. Following a series of lemmas in Appendix A, the proof of Theorem 1 is given in Appendix A.

A by-product of the proof for Theorem 1, given in Appendix A, shows that $\boldsymbol{\mathcal{I}}(\boldsymbol{\beta}_0) = \lim_{n \rightarrow \infty} N_n^{-1} \mathbf{X}^\top (\mathbf{I} - \mathbf{S})^\top \boldsymbol{\Gamma}^{-1} (\mathbf{I} - \mathbf{S}) \mathbf{X}$. Thus, we use $N_n^{-1} \mathbf{X}^\top (\mathbf{I} - \mathbf{S})^\top \boldsymbol{\Gamma}^{-1} (\mathbf{I} - \mathbf{S}) \mathbf{X}$ as a finite sample approximation of $\boldsymbol{\Pi}$, the asymptotic variance of ${}^n\hat{\boldsymbol{\beta}}$. In contrast, for ${}^n\hat{\boldsymbol{\theta}}$, it can be shown that the asymptotic variance is the same as that for the case when the temporal function $f(\cdot)$ is assumed to be known.

Further, recall that $\hat{\mathbf{F}}(t) = \boldsymbol{\omega}(t)(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})$ is the estimate of $\mathbf{F}(t) = (f(t), hf'(t))^\top$, where $\boldsymbol{\omega}(t) = (\mathbf{D}_t^\top \mathbf{K}_t \mathbf{D}_t)^{-1} \mathbf{D}_t^\top \mathbf{K}_t$. The following Theorem 2 establishes the asymptotic normality of $\hat{\mathbf{F}}(t)$. The proof of Theorem 2 is given in Appendix A.

Theorem 2. Suppose $f^{(3)}(t)$ is bounded. Under (C.1)–(C.13) in the Appendix A, we have, as $n \rightarrow \infty$,

$$(N_n h)^{1/2} \left\{ \hat{\mathbf{F}}(t) - \mathbf{F}(t) - (1/2)h^2 \begin{pmatrix} \mu_2 f''(t) \\ 0 \end{pmatrix} + o(h^2) \right\} \xrightarrow{D} \mathcal{N} \left(\mathbf{0}, \begin{pmatrix} 1 & 0 \\ 0 & \mu_2^{-1} \end{pmatrix} \Delta_t \begin{pmatrix} 1 & 0 \\ 0 & \mu_2^{-1} \end{pmatrix} \right),$$

for $t \in (0, 1)$, where $\mu_k = \int_{-\infty}^{\infty} x^k K(x) dx$.

Similar to [Theorems 1 and 2](#), we can show that, when β is known, the asymptotic properties of $\hat{\theta}$ and $\hat{F}(t)$ remain the same as in [Theorem 2](#). This may be expected, since $\hat{\beta}$ is root- N_n consistent. Further, Δ_t is generally unknown in practice. As suggested by [\(C.9\)](#), Δ_t can be approximated by $N_n^{-1} h q(t)^{-2} \mathbf{k}_t^\top \Gamma \mathbf{k}_t$, where $\mathbf{k}_t = \{K_h(t_i - t) \{(t_i - t)/h\}^{j-1}\}_{i,j=1}^{N_n, 2}$ is an $N_n \times 2$ matrix. Since $q(t)$ represents the density of the sampling time points, we may use a kernel density method to estimate $q(t)$. An alternative is to estimate $q(t)$ by $v_{0,t}/N_n$, where $v_{0,t} = \sum_{i=1}^{N_n} K_h(t_i - t)$. [Lemma 1](#) in [Appendix A](#) shows that such an approximation is reasonable. In the remainder of this paper, we will refer to the former approximation as a kernel density approximation and the latter a plug-in approximation.

4. Selection of kernel and bandwidth

4.1. Theoretically optimal bandwidth

The selection of bandwidth is crucial in kernel smoothing and thus we derive a theoretically optimal bandwidth. By the results in [Theorem 2](#), the asymptotic mean squared error (AMSE) of $\hat{f}(t)$ is

$$\text{AMSE}(t) = (1/4)h^4 \mu_2^2 f''(t)^2 + (N_n h)^{-1} (1, 0) \Delta_t (1, 0)^\top,$$

and the asymptotic weighted mean integrated squared error is

$$\text{AMISE}(h) = \int_0^1 \text{AMSE}(t) q(t) dt = (1/4)h^4 \mu_2^2 \int_0^1 f''(t)^2 q(t) dt + (N_n h)^{-1} \int_0^1 (1, 0) \Delta_t (1, 0)^\top q(t) dt.$$

Viewing the density function $q(t)$ as a weight function, we obtain an asymptotically optimal bandwidth as

$$h_{\text{opt}} = N_n^{-1/5} \mu_2^{-2/5} \left\{ \frac{\int_0^1 (1, 0) \Delta_t (1, 0)^\top q(t) dt}{\int_0^1 f''(t)^2 q(t) dt} \right\}^{1/5}, \quad (7)$$

where the convergence rate is $N_n^{2/5}$ and is the nonparametric optimal rate [\[40\]](#).

The asymptotically optimal bandwidth h_{opt} above depends on several unknown quantities: Δ_t in the asymptotic variance of $\hat{F}(t)$, the density of sampling time points $q(t)$, and the second-order derivative of the temporal function $f''(t)$; thus, it is not straightforward to estimate h_{opt} . When $\Gamma = \sigma^2 \mathbf{I}$ (i.e., the process assumes spatio-temporal independence), Δ_t can be expressed as $\sigma^2 q(t)^{-1} \text{diag}\{\int_{-\infty}^{\infty} K(u)^2 du, \int_{-\infty}^{\infty} u^2 K(u)^2 du\}$. A rule of thumb for bandwidth selection in this case is available [see, e.g., [13](#)]. The idea is to plug in the estimates of σ^2 and $f''(t)$ to obtain an approximation of h_{opt} . Specifically, after a pilot global polynomial regression of degree 4 is fitted, σ^2 is estimated by the standardized residual sum of squares, and the estimate of $f''(t)$ is obtained by differentiating the resulting global fit. However, for a spatio-temporally correlated error process, the covariance matrix Γ needs to be estimated, and this rule of thumb is not directly applicable. Hence, a more practical bandwidth selection procedure is needed.

4.2. Practical bandwidth selection

Under model [\(1\)](#), we have

$$y^*(\mathbf{s}, t) = f(t) + \varepsilon_1(\mathbf{s}, t) + \varepsilon_2(\mathbf{s}, t), \quad \mathbf{s} \in \mathcal{R}, \quad t \in \mathcal{T}. \quad (8)$$

We use a leave-one-out cross-validation criterion [CV; [44](#)]. A straightforward calculation reveals that $\text{CV}(h) = N_n^{-1} \sum_{i=1}^{N_n} \left\{ \frac{y_i^* - \hat{f}(t_i)}{1 - S_{ii}} \right\}^2$, where S_{ii} is the (i, i) th element of the smoother matrix \mathbf{S} . However, as will be seen in the following theorem, cross-validation is asymptotically biased in the presence of correlated errors for most commonly used kernels with $K(0) \neq 0$.

Theorem 3. Under Assumptions [\(C.1\)–\(C.13\)](#) in the [Appendix A](#), if there exists a sequence $C_n > 0$ such that $C_n h^{-1} \rightarrow 0$ and $1/(B_n \zeta_n) \int_{B_n C_n}^{\infty} \gamma_1(u) du \rightarrow 0$, as $n \rightarrow \infty$, then we have

$$\text{E}\{\text{CV}(h)\} = N_n^{-1} \sum_{i=1}^{N_n} \text{E}\{f(t_i) - \hat{f}^{(-i)}(t_i)\}^2 + \overline{\sigma^2} - K(0) \left\{ (2/N_n) \sum_{i=1}^{N_n} \sum_{\substack{j \neq i \\ |t_j - t_i| < C_n}} \frac{\text{Cov}(\varepsilon_i, \varepsilon_j)}{b(t_i) - K(0)} \right\} + o(1/(N_n h)),$$

where $\hat{f}^{(-i)}(t_i)$ is the leave-one-out estimator with the i th observation deleted for estimation, $\overline{\sigma^2} = N_n^{-1} \sum_{i=1}^{N_n} \text{Var}(Y_i)$, $\text{CV}(h) = N_n^{-1} \sum_{i=1}^{N_n} \{y_i^* - \hat{f}^{(-i)}(t_i)\}^2$ and $b(t_i) = N_n q(t_i) h (\mu_{0,t_i} \mu_{2,t_i} - \mu_{1,t_i}^2) \mu_{2,t_i}^{-1}$.

The proof of [Theorem 3](#) is given in [Appendix A](#). [Theorem 3](#) provides a theoretical basis for the choice of kernel functions. In practice, we propose the following procedure for the selection of bandwidth h .

- (i) For a predetermined bandwidth h_0 and a kernel function K_{h_0} , obtain the estimated regression coefficients $\tilde{\beta}$ by the profile likelihood method.
- (ii) For given a kernel function, find the bandwidth h_{opt} that minimizes the cross-validation criterion

$$CV(h) = N_n^{-1} \sum_{i=1}^{N_n} \left(\frac{\tilde{y}_i^* - \tilde{f}_i^*}{1 - S_{ii}} \right)^2, \quad (9)$$

where $\tilde{y}_i^* = y_i - \mathbf{x}(\mathbf{s}_i, t_i)^\top \tilde{\beta}$, and \tilde{f}_i^* is the profile likelihood estimate of (8).

- (iii) Use h_{opt} and the kernel function from Step (ii) to obtain the desired estimates of both the regression coefficients β and the covariance function parameters θ .

As to be illustrated in a simulation study, the estimate of β is not very sensitive to the choices of bandwidth and kernel function in Step (i). Thus, we suggest to use a pilot bandwidth to yield an underestimate of f and consequently an estimate of β . In Steps (ii) and (iii), we use a bimodal kernel $K_2(u) = 2\pi^{-1/2}u^2 \exp(-u^2)$; see Fig. B of the supplementary material [12]. Unlike the more commonly used kernels (e.g., Gaussian or Epanechnikov kernel), the bimodal kernel satisfies $K_2(0) = 0$, which can mitigate the influence of the spatio-temporal correlation.

A popular alternative to the cross-validation criterion (9) is the generalized cross-validation [GCV;20] criterion, in which S_{ii} is replaced by $N_n^{-1} \text{tr}(\mathbf{S})$. For dependent data, Francisco-Fernandez and Opsomer [15] proposed a bias-corrected generalized cross-validation criterion (GCV_c), replacing S_{ii} by $N_n^{-1} \text{tr}(\mathbf{S}\mathbf{R}(\theta))$; that is,

$$\text{GCV}_c(h) = \frac{\sum_{i=1}^{N_n} (\tilde{y}_i^* - \tilde{f}_i^*)^2}{N_n \{1 - N_n^{-1} \text{tr}(\mathbf{S}\mathbf{R}(\theta))\}^2}, \quad (10)$$

where $\mathbf{R}(\theta)$ is a correlation matrix. In practice, a pilot estimate of the covariance parameter vector is required; however, the choice of such an estimate is not obvious, and would impact the overall estimation performance. To ensure the performance of parameter estimation in covariance function, for each candidate bandwidth h , we compute the corresponding estimate of θ and obtain an estimated GCV_c criterion, denoted by GCV_{ce} . As further demonstrated in the simulation study, the results based on the cross-validation and GCV_{ce} are similar, although GCV_{ce} is computationally more expensive.

5. Simulation study

5.1. Simulation set-up

We sample N_s locations uniformly from the spatial domain $[0, 1]^2$, where $N_s \in \{20, 40, 60\}$. At each sampling location, we randomly sample 4% from the grid of time points $(i - 1/2)/1000$, $i \in \{1, \dots, 1000\}$. The selected locations and time are labeled as $\{(\mathbf{s}_1, t_1), \dots, (\mathbf{s}_{N_n}, t_{N_n})\}$. For $N_s \in \{20, 40, 60\}$, the sample sizes are $N_n \in \{806, 1644, 2449\}$, respectively. The space-time coordinates will remain fixed across iterations once generated.

For the regression mean function and the semiparametric mean function, the vector of regression coefficients is $\beta = (4, 3, 2, 1)^\top$. The covariates are drawn (once) from a multivariate normal distribution with zero mean, unit variance, and a cross-covariate correlation of 0.5. Each covariate is standardized to have zero sample mean and unit sample variance. Further, the nonparametric temporal function in the semiparametric mean function is $f(t) = 2\{1 - \cos(2\pi t)\}$.

We then draw a realization from the mean zero Gaussian error process $\varepsilon(\mathbf{s}, t)$ using three different covariance functions. The first covariance function is an exponential spatio-temporal covariance function

$$\text{Cov}\{\varepsilon(\mathbf{s}_i, t_i), \varepsilon(\mathbf{s}_j, t_j)\} = \begin{cases} \sigma^2(1 - c) \exp\{-\varrho_{1,n}\|\mathbf{s}_i - \mathbf{s}_j\|/c_s - \varrho_{2,n}|t_i - t_j|/c_t\}, & \text{if } i \neq j, \\ \sigma^2, & \text{if } i = j. \end{cases}$$

Here, σ^2 is the variance of $\varepsilon(\mathbf{s}, t)$, $c \in [0, 1]$ is the proportion of random noise such that $c\sigma^2$ is the nugget effect, and c_s and c_t are the positive spatial and temporal range parameters, respectively. We take $\sigma^2 = 9.0$, $c = 0.2$, $c_s = 1$ and $c_t = 1$. This covariance function is stationary and separable and we denote it as COV-1.

Next, we consider a generalized spatio-temporal Matérn covariance function [5]:

$$\gamma_n((\mathbf{s}, t), (\mathbf{s}', t'); \theta) = \begin{cases} \frac{D(\mathbf{s}, t)D(\mathbf{s}', t')\sigma^2\theta_3^{d/2}2^{1-\nu}}{(\theta_1^2u_2^2+1)^\nu(\theta_1^2u_2^2+\theta_3)^{d/2}\Gamma(\nu)} m(\mathbf{u}_1, u_2)^\nu K_\nu\{m(\mathbf{u}_1, u_2)\}, & \text{if } \|\mathbf{u}_1\| > 0, \\ \frac{D(\mathbf{s}, t)D(\mathbf{s}', t')\sigma^2\theta_3^{d/2}}{(\theta_1^2u_2^2+1)^\nu(\theta_1^2u_2^2+\theta_3)^{d/2}}, & \text{if } \|\mathbf{u}_1\| = 0, |u_2| > 0, \\ D(\mathbf{s}, t)^2\sigma^2 + c\sigma^2, & \text{if } \|\mathbf{u}_1\| = 0, |u_2| = 0, \end{cases} \quad (11)$$

where $\mathbf{u}_1 = \rho_{1,n}(\mathbf{s}' - \mathbf{s})$ and $u_2 = \rho_{2,n}(t' - t)$. In this covariance function, $m(\mathbf{u}_1, u_2) = \theta_2\|\mathbf{u}_1\|\{(\theta_1^2u_2^2+1)/(\theta_1^2u_2^2+\theta_3)\}^{1/2}$ and $K_\nu(\cdot)$ is the modified Bessel function of the second kind of order ν . Here, θ_1 and θ_2 are nonnegative range parameters of time and space respectively, $\theta_3 > 0$ is a separability parameter. The point-wise variance of $\varepsilon(\mathbf{s}, t)$ is $D(\mathbf{s}, t)D(\mathbf{s}', t')\sigma^2 + c\sigma^2$,

where $c\sigma^2$ accounts for the nugget effect. The parameter ν in $K_\nu(\cdot)$ controls the smoothness of the covariance. If we let $\nu = 1/2$ and $\theta_3 = 1$, then (11) reduces to

$$\text{Cov}\{\varepsilon(\mathbf{s}_i, t_i), \varepsilon(\mathbf{s}_j, t_j)\} = \begin{cases} D(\mathbf{s}_i, t_i)D(\mathbf{s}_j, t_j) \frac{\sigma^2}{\{a^2|\varrho_{2,n}(t_i-t_j)|^2+1\}^{3/2}} \exp\{-b\varrho_{1,n}\|\mathbf{s}_i - \mathbf{s}_j\|\}, & \text{if } i \neq j; \\ D(\mathbf{s}_i, t_i)D(\mathbf{s}_j, t_j)\sigma^2 + c\sigma^2, & \text{if } i = j. \end{cases} \quad (12)$$

Similarly, σ^2 is the variance of $\varepsilon(\mathbf{s}, t)$, $c \in [0, 1]$ is the proportion of random noise such that $c\sigma^2$ is the nugget effect, and a and b are the positive temporal and spatial range parameters, respectively. Here, $D(\mathbf{s}_i, t_i) = dt_i + 1$ varies by time, resulting in a nonstationarity covariance function. We set $\sigma^2 = 9$, $c = 0.2$, $a = 1$, $b = 1$ and $d = 1$. This covariance function is still separable in space and time, which we refer to as COV-2.

The third covariance function we considered is a slight modification of COV-2, with $D(\mathbf{s}_i, t_i) = dt_i + es_{1i} + fs_{2i} + 1$. We set $\sigma^2 = 9.0$, $c = 0.2$, $a = 1$, $b = 1$, $d = 0.5$, $e = 0.5$ and $f = 0.5$. This covariance function is nonstationary, nonseparable and asymmetric, referred as COV-3.

For each combination of the sample size and the covariance function, we generate 400 simulation replicates. We also consider a special case of the semiparametric mean function where the temporal function f is assumed to be zero. As will be demonstrated later, this case will serve as a benchmark in the comparison of the estimation for β and θ .

For each simulated data set, a predetermined bandwidth $h_0 = 0.05$ is used to obtain an initial estimate of β . The estimate of the optimal bandwidth \hat{h} is then determined by minimizing the cross-validation criterion (9) over a predetermined grid of bandwidth values. Given the estimated optimal bandwidth, the profile likelihood estimates $\hat{\beta}$, $\hat{\theta}$ and $\hat{f}(\cdot)$ are obtained. We further consider two variants of the GCV for determining the bandwidth in (9): GCV_c and GCV_{ce} as described in Section 4.

The profile likelihood method (PLE) results are compared with two alternative methods, namely, ALT₁ and ALT₂. In ALT₁, the parameter estimates and the estimate of the temporal function are obtained by the profile likelihood method ignoring the spatio-temporal dependence. In ALT₂, the regression coefficients β and the covariance parameters θ are estimated by the classical maximum likelihood method assuming the temporal trend $f(\cdot)$ is known. That is, ALT₂ is essentially the maximum likelihood method under the model with the regression mean function.

To assess the performance of estimation by the different methods under the different bandwidth selection criteria, we compute the means and the standard deviations (SD) of $\hat{\beta}$ and $\hat{\theta}$ from the 400 simulated data sets. We also compute the estimated standard errors of the parameters for each simulated data set based on the information matrix in Theorem 1 and report the mean estimated standard errors (SDm). For ALT₁, we use $\Gamma(\hat{\theta}) = \hat{\sigma}^2 \mathbf{I}$ to calculate SDm. In addition, for the estimated temporal function \hat{f} , we calculate the average squared error (ASE) for each simulated data set, defined as

$$\text{ASE} = N_{\text{grid}}^{-1} \sum_{i=1}^{N_{\text{grid}}} \{f(t_{i,\text{grid}}) - \hat{f}(t_{i,\text{grid}})\}^2,$$

where $t_{i,\text{grid}} = (i - 1/2)/N_{\text{grid}}$ for $i \in \{1, \dots, N_{\text{grid}}\}$ and $N_{\text{grid}} = 1000$.

Finally, we generate an additional 10% new sampling locations and, at each new sampling location, new sampling time points are generated as in the simulation set-up. At these new sampling locations and time points, new observations denoted as $y_{i,\text{new}}$ are generated and let $\tilde{y}_{i,\text{new}}$ denote the predicted value at the i th new sampling location and time, where $i \in \{1, \dots, N_{\text{new}}\}$, and N_{new} is the total number of new sampling locations and time points. We use the mean squared prediction error (MSPE) to evaluate the performance of the various methods as

$$\text{MSPE} = N_{\text{new}}^{-1} \sum_{i=1}^{N_{\text{new}}} (y_{i,\text{new}} - \tilde{y}_{i,\text{new}})^2,$$

The results are provided in Tables 1–3, the last two rows of which give the average values of ASE and MSPE.

5.2. Simulation results

As shown in Table 1 for the first scenario of the spatio-temporal covariance function (COV-1), the bandwidths chosen by the three selection criteria, CV, GCV_c and GCV_{ce} , are similar for the profile likelihood method. For parameter estimation, both the accuracy and the precision increase as the sample size increases. The empirical standard deviations are well approximated by the standard errors, supporting the information-based asymptotic variance in Theorem 1. Further, under different bandwidth selection criteria, similar ASE and MSPE values are obtained, which may not be surprising due to the similar choices of bandwidths and hence similar estimates.

For the estimation of the regression coefficients, our method PLE and the two alternative methods ALT₁ and ALT₂ have comparable estimation bias, which suggests that the accuracy of $\hat{\beta}$ is not sensitive to the assumption of covariance structure. However, the simulation standard deviations from ALT₁ are larger than those from PLE and ALT₂, indicating noticeable gain of statistical efficiency in the parameter estimation by accounting for spatio-temporal dependence. In addition, ALT₁ has much larger MSPE and thus poorer prediction than PLE and ALT₂. For estimating the temporal function, the ASEs for ALT₁ and PLE are similar; both decrease as the sample size increases.

Table 1

Sample mean, sample standard deviation (SD), averaged information matrix based standard deviation (SD_m) of regression and covariance parameters, averaged squared error (ASE) and mean-squared prediction error (MSPE) for three sample sizes with $N_s = 20, 40, 60$ using bimodal kernel for COV-1 and for three bandwidth selection criteria, CV, GCV_c, GCV_{ce} under profile likelihood estimation (PLE). Comparison is made with two alternatives ALT₁ and ALT₂.

Term	Truth	$N_s = 20$						$N_s = 40$						$N_s = 60$					
		PLE			ALT ₁		ALT ₂			PLE			ALT ₁		ALT ₂				
		CV	GCV _c	GCV _{ce}	CV	–	CV	GCV _c	GCV _{ce}	CV	–	CV	GCV _c	GCV _{ce}	CV	–	CV	GCV _c	GCV _{ce}
h	–	0.079	0.081	0.082	0.079	–	0.072	0.073	0.073	0.072	–	0.068	0.069	0.069	0.068	–			
β_1	4.0	3.985	3.985	3.985	3.991	3.985	4.002	4.002	4.002	3.993	4.002	4.002	4.002	4.002	4.003	4.002			
SD		0.104	0.104	0.104	0.120	0.104	0.080	0.080	0.080	0.104	0.080	0.067	0.067	0.067	0.082	0.067			
SD _m		0.112	0.112	0.112	0.132	0.112	0.080	0.080	0.080	0.095	0.080	0.066	0.066	0.066	0.078	0.066			
β_2	3.0	3.017	3.017	3.017	3.017	3.018	3.009	3.009	3.009	3.004	3.010	2.996	2.996	2.996	2.989	2.996			
SD		0.123	0.123	0.123	0.140	0.122	0.075	0.075	0.075	0.095	0.075	0.065	0.065	0.065	0.079	0.065			
SD _m		0.116	0.116	0.116	0.136	0.117	0.078	0.078	0.078	0.094	0.078	0.065	0.065	0.065	0.077	0.065			
β_3	2.0	2.006	2.006	2.006	2.003	2.005	1.986	1.986	1.986	1.982	1.986	1.997	1.997	1.997	2.001	1.997			
SD		0.109	0.109	0.109	0.127	0.109	0.074	0.074	0.074	0.087	0.074	0.067	0.067	0.067	0.082	0.067			
SD _m		0.114	0.114	0.114	0.132	0.114	0.081	0.081	0.081	0.096	0.081	0.064	0.064	0.064	0.077	0.064			
β_4	1.0	0.996	0.996	0.995	0.988	0.996	1.002	1.002	1.002	1.010	1.002	1.000	1.000	1.000	0.999	1.000			
SD		0.110	0.110	0.110	0.131	0.110	0.077	0.077	0.077	0.097	0.077	0.065	0.065	0.065	0.082	0.065			
SD _m		0.115	0.115	0.115	0.134	0.115	0.079	0.079	0.079	0.094	0.079	0.065	0.065	0.065	0.077	0.065			
σ^2	9.0	9.111	9.112	9.112	9.101	8.944	9.120	9.120	9.120	9.113	8.997	9.055	9.055	9.055	9.054	8.963			
SD		0.569	0.569	0.569	0.589	0.525	0.387	0.387	0.387	0.394	0.372	0.325	0.324	0.324	0.334	0.317			
SD _m		0.546	0.547	0.547	–	0.528	0.395	0.395	0.395	–	0.384	0.319	0.319	0.319	–	0.313			
c	0.2	0.209	0.209	0.209	–	0.202	0.201	0.201	0.201	–	0.196	0.197	0.197	0.197	–	0.193			
SD		0.077	0.077	0.077	–	0.078	0.047	0.047	0.048	–	0.047	0.043	0.043	0.043	–	0.043			
SD _m		0.074	0.074	0.074	–	0.077	0.047	0.047	0.047	–	0.048	0.040	0.040	0.040	–	0.041			
c_s	1.0	1.090	1.090	1.091	–	1.025	1.052	1.052	1.052	–	1.007	1.027	1.027	1.027	–	0.994			
SD		0.224	0.224	0.224	–	0.198	0.134	0.134	0.134	–	0.124	0.103	0.103	0.103	–	0.098			
SD _m		0.213	0.213	0.213	–	0.200	0.127	0.127	0.127	–	0.121	0.099	0.099	0.099	–	0.095			
c_t	1.0	1.087	1.088	1.088	–	1.029	1.047	1.048	1.048	–	1.008	1.016	1.016	1.016	–	0.987			
SD		0.243	0.244	0.245	–	0.224	0.142	0.142	0.142	–	0.134	0.112	0.112	0.112	–	0.108			
SD _m		0.213	0.213	0.213	–	0.204	0.139	0.139	0.139	–	0.134	0.112	0.112	0.112	–	0.109			
ASE	–	0.266	0.265	0.266	0.265	–	0.183	0.183	0.183	0.182	–	0.149	0.149	0.150	0.149	–			
MSPE	–	6.501	6.501	6.501	9.245	6.484	7.403	7.403	7.403	9.169	7.386	6.934	6.933	6.933	9.090	6.918			

When the sample size is smaller, PLE has less accuracy and precision in the estimation than ALT₂. In particular, both the standard deviations and the standard errors of the estimates from PLE are considerably larger than those of ALT₂, for all the covariance parameters except the nugget proportion c . When the sample size is larger, PLE and ALT₂ have similar estimation results. In particular, the standard deviations and the standard errors of the estimates from PLE are similar to ALT₂, supporting that the asymptotic variance of $\hat{\theta}$ under the semiparametric mean function is the same as the regression mean function, as shown in Theorem 1. Moreover, ALT₂ has slightly better prediction than PLE due to possible bias in the estimation of $f(\cdot)$ in PLE.

Tables 2 and 3 show results for the second and the third scenario of the spatio-temporal covariance function, COV-2 and COV-3, respectively. Similar conclusions can be drawn. Particularly, in the presence of non-separability and nonstationarity in the spatio-temporal covariance function, the finite-sample performance of the estimation for the semiparametric mean function is sound and supports the asymptotic results. The bandwidths selected by the three criteria, CV, GCV_c and GCV_{ce}, are very similar and so are the resulting estimates. Unlike COV-1 and COV-2, the prediction under COV-3 changes greatly for different sample sizes, which may be attributed to the nonstationarity in space with very different variances at different new spatial locations where the observations are predicted. Tables D–E in Section 3 of the Supplementary Material show that the regression coefficient estimates are robust against the choice of the kernel and the initial bandwidth. For a predetermined bandwidth, it can be seen that, different kernel functions in Step (i) of the bandwidth selection procedure yield very similar results. Moreover, those results are similar to the benchmark case when $\beta = \beta$.

As demonstrated in Theorem 3, bimodal kernels can effectively alleviate the influence of correlated errors on bandwidth selection. To see this, we compare the results from a bimodal kernel with those from a Gaussian kernel. The estimation results under the Gaussian kernel for the first scenario of spatio-temporal covariance function (COV-1) are given in Table F of the Supplementary Material. Unlike the bimodal kernel, the bandwidths selected by the three criteria, CV, GCV_c and GCV_{ce}, can be quite different. In particular, CV selects much smaller bandwidth than GCV_c and GCV_{ce}, supporting the fact that cross-validation does not handle correlation well for most commonly-used kernels with $K(0) \neq 0$. The regression coefficient estimates are similar for all three bandwidth selection criteria, which suggests that the estimation of β is not sensitive to the choice of bandwidth. However, the estimates of covariance parameters are greatly affected by the bias in the bandwidth selection, with CV having the largest bias in parameter estimation and the largest ASE in the estimation of the temporal function, among the three criteria. As an alternative of CV in the presence of spatio-temporal correlation, GCV_c produces a much larger bandwidth, although the resulting estimates are not as accurate as those from the bimodal kernel function.

Table 2

Sample mean, sample standard deviation (SD), averaged information matrix based standard deviation (SD_m) of regression and covariance parameters, averaged squared error (ASE), and mean-squared prediction error (MSPE) for three sample sizes with $N_s = 20, 40, 60$ using bimodal kernel for COV-2 and for three bandwidth selection criteria, CV, GCV_c, GCV_{ce} under profile likelihood estimation (PLE). Comparison is made with two alternatives ALT₁ and ALT₂.

Term	Truth	$N_s = 20$						$N_s = 40$						$N_s = 60$					
		PLE			ALT ₁		ALT ₂	PLE			ALT ₁		ALT ₂	PLE			ALT ₁		ALT ₂
		CV	GCV _c	GCV _{ce}	CV	–		CV	GCV _c	GCV _{ce}	CV	–		CV	GCV _c	GCV _{ce}	CV	–	
h	–	0.100	0.103	0.104	0.100	–	–	0.089	0.091	0.091	0.089	–	–	0.083	0.086	0.086	0.083	–	–
β_1	4.0	3.969	3.969	3.969	3.987	3.969	–	4.002	4.002	4.002	3.987	4.003	–	4.004	4.004	4.004	4.005	4.004	–
SD		0.136	0.137	0.137	0.197	0.137	–	0.102	0.102	0.102	0.169	0.102	–	0.089	0.089	0.089	0.129	0.089	–
SD _m		0.146	0.146	0.146	0.209	0.146	–	0.104	0.104	0.104	0.151	0.104	–	0.086	0.086	0.086	0.123	0.086	–
β_2	3.0	3.020	3.020	3.020	3.027	3.023	–	3.010	3.010	3.010	3.008	3.012	–	2.994	2.994	2.994	2.983	2.994	–
SD		0.165	0.165	0.165	0.218	0.164	–	0.097	0.097	0.097	0.155	0.097	–	0.083	0.083	0.083	0.128	0.082	–
SD _m		0.153	0.153	0.153	0.215	0.153	–	0.101	0.101	0.101	0.150	0.101	–	0.082	0.082	0.082	0.122	0.082	–
β_3	2.0	2.011	2.011	2.011	1.998	2.008	–	1.984	1.984	1.984	1.970	1.984	–	1.993	1.993	1.993	2.000	1.993	–
SD		0.146	0.146	0.146	0.201	0.146	–	0.100	0.100	0.100	0.136	0.100	–	0.088	0.088	0.088	0.130	0.088	–
SD _m		0.152	0.152	0.152	0.209	0.152	–	0.106	0.106	0.106	0.154	0.106	–	0.082	0.082	0.082	0.121	0.082	–
β_4	1.0	1.003	1.003	1.003	0.983	1.004	–	1.002	1.002	1.002	1.021	1.002	–	1.002	1.002	1.002	1.000	1.002	–
SD		0.144	0.144	0.144	0.208	0.144	–	0.101	0.101	0.101	0.161	0.102	–	0.085	0.085	0.085	0.130	0.085	–
SD _m		0.148	0.149	0.149	0.213	0.148	–	0.102	0.102	0.102	0.150	0.102	–	0.084	0.084	0.084	0.122	0.084	–
σ^2	9.0	9.075	9.073	9.071	22.751	8.934	–	9.198	9.199	9.199	23.316	9.094	–	9.163	9.159	9.159	22.736	9.085	–
SD		1.610	1.609	1.608	1.658	1.543	–	1.131	1.133	1.133	1.217	1.114	–	0.871	0.869	0.869	0.995	0.864	–
SD _m		1.569	1.569	1.569	–	1.544	–	1.114	1.114	1.114	–	1.100	–	0.887	0.887	0.887	–	0.879	–
c	0.2	0.229	0.229	0.230	–	0.227	–	0.200	0.200	0.200	–	0.199	–	0.195	0.195	0.195	–	0.194	–
SD		0.121	0.121	0.121	–	0.120	–	0.060	0.060	0.060	–	0.061	–	0.051	0.051	0.051	–	0.051	–
SD _m		0.102	0.102	0.102	–	0.103	–	0.061	0.061	0.061	–	0.062	–	0.049	0.049	0.049	–	0.050	–
a	1.0	0.980	0.980	0.980	–	0.996	–	0.991	0.991	0.991	–	1.002	–	1.002	1.002	1.002	–	1.010	–
SD		0.117	0.117	0.117	–	0.118	–	0.070	0.070	0.070	–	0.071	–	0.058	0.058	0.058	–	0.058	–
SD _m		0.101	0.101	0.101	–	0.104	–	0.069	0.069	0.069	–	0.071	–	0.058	0.058	0.058	–	0.059	–
b	1.0	0.973	0.973	0.973	–	1.001	–	0.984	0.984	0.984	–	1.006	–	0.992	0.992	0.992	–	1.009	–
SD		0.139	0.139	0.139	–	0.140	–	0.089	0.089	0.089	–	0.089	–	0.070	0.070	0.070	–	0.070	–
SD _m		0.131	0.131	0.131	–	0.135	–	0.086	0.086	0.086	–	0.087	–	0.069	0.069	0.069	–	0.070	–
d	1.0	1.020	1.020	1.020	–	1.021	–	0.999	0.999	0.999	–	1.000	–	0.991	0.992	0.992	–	0.992	–
SD		0.246	0.247	0.247	–	0.242	–	0.167	0.167	0.167	–	0.167	–	0.127	0.127	0.127	–	0.127	–
SD _m		0.238	0.238	0.238	–	0.237	–	0.162	0.162	0.162	–	0.161	–	0.129	0.129	0.129	–	0.129	–
ASE	–	0.663	0.659	0.657	0.655	–	–	0.451	0.451	0.451	0.450	–	–	0.365	0.366	0.367	0.366	–	–
MSPE	–	13.454	13.456	13.456	24.114	13.429	–	15.646	15.645	15.646	21.682	15.609	–	15.629	15.629	15.629	24.123	15.598	–

Finally, in Tables 1–3, the standard deviations and the standard errors of $\hat{\beta}$ from PLE are similar to those from ALT₂. This is as expected, since the design matrix in our setting does not vary by time. In Section 1.2 of the Supplementary Material, we investigate a design matrix that varies over time. From Table G in the Supplementary Material, it can be seen that the standard deviations of $\hat{\beta}$ from PLE are larger than those from ALT₂, which indicates a loss of statistical efficiency in the estimation of β when the unknown temporal function is estimated. This finding is consistent with the standard error formula in Theorem 1. In addition, we consider a nonseparable but stationary covariance function. The simulation results are provided in Section 1.3 of the Supplementary Material, and similar lessons can be learned.

6. Data example

To illustrate our methodology, we consider a data set collected by static sensors at fixed sampling locations in time and roving sensors traversing the spatial domain in time in an engine facility for evaluating the intensity level of noise as an occupational hazard [25,31]. We focus on the observations between 10:29:00 am and 11:24:00 am when all the sensors are operating. As shown in Fig. 1, there are 56 observations, one per minute, for each of the 17 static sensors. For the two roving sensors, there are a total of 179 observations, observed at irregular time points. Therefore, the total sample size is $N_n = 1131$.

We consider the semiparametric mean function (12) with the generalized spatio-temporal Matérn error covariance function (12). More specifically, for $\mathbf{s} = (s_1, s_2)$, we have

$$y(\mathbf{s}, t) = \beta_1 s_1 + \beta_2 s_2 + f(t) + \varepsilon(\mathbf{s}, t),$$

where the regression is on the coordinates of the spatial location \mathbf{s} , the temporal function f is nonparametric, and the zero-mean error process $\varepsilon(\mathbf{s}, t)$ has the spatio-temporal covariance function (12). We fit three spatio-temporal covariance functions: $D_1(\mathbf{s}, t) = 1$ for stationarity, $D_2(\mathbf{s}, t) = 1 + dt$ and $D_3(\mathbf{s}, t) = 1 + dt + e(t - \kappa)_+$ for nonstationary. In the latter two nonstationary cases, for any fixed time point t_0 , $\varepsilon(\mathbf{s}, t_0)$ is spatially stationary. For $D_3(\mathbf{s}, t)$, κ is chosen around 11:02:00 am, which is expected to capture the temporal change due to an engine shutdown.

Table 3
Sample mean, sample standard deviation (SD), averaged information matrix based standard deviation (SD_m) of regression and covariance parameters, averaged squared error (ASE), and mean-squared prediction error (MSPE) for three sample sizes with $N_s = 20, 40, 60$ using bimodal kernel for COV-3 and for three bandwidth selection criteria, CV, GCV_c, GCV_{ce} under profile likelihood estimation (PLE). Comparison is made with two alternatives ALT₁ and ALT₂.

Term	Truth	$N_s = 20$						$N_s = 40$						$N_s = 60$					
		PLE			ALT ₁		ALT ₂	PLE			ALT ₁		ALT ₂	PLE			ALT ₁		ALT ₂
		CV	GCV _c	GCV _{ce}	CV	–		CV	GCV _c	GCV _{ce}	CV	–		CV	GCV _c	GCV _{ce}	CV	–	
h	–	0.106	0.108	0.109	0.106	–	–	0.095	0.098	0.099	0.095	–	–	0.087	0.090	0.090	0.087	–	–
β_1	4.0	3.964	3.963	3.963	3.984	3.963	–	4.003	4.003	4.003	3.984	4.003	–	4.005	4.005	4.005	4.004	4.005	–
SD		0.154	0.154	0.154	0.220	0.154	–	0.114	0.114	0.114	0.193	0.114	–	0.098	0.098	0.098	0.148	0.098	–
SD _m		0.165	0.165	0.165	0.238	0.165	–	0.116	0.116	0.116	0.173	0.116	–	0.096	0.096	0.096	0.139	0.096	–
β_2	3.0	3.025	3.025	3.025	3.029	3.027	–	3.014	3.014	3.014	3.011	3.015	–	2.992	2.992	2.992	2.983	2.992	–
SD		0.184	0.184	0.184	0.248	0.183	–	0.107	0.107	0.107	0.181	0.107	–	0.092	0.092	0.092	0.144	0.092	–
SD _m		0.171	0.171	0.171	0.246	0.172	–	0.112	0.112	0.112	0.171	0.112	–	0.092	0.092	0.092	0.138	0.092	–
β_3	2.0	2.014	2.014	2.014	1.994	2.011	–	1.981	1.981	1.981	1.968	1.981	–	1.991	1.991	1.991	2.002	1.991	–
SD		0.162	0.162	0.162	0.230	0.161	–	0.111	0.111	0.111	0.158	0.111	–	0.097	0.097	0.097	0.148	0.097	–
SD _m		0.172	0.172	0.172	0.238	0.172	–	0.117	0.117	0.117	0.176	0.117	–	0.092	0.092	0.092	0.137	0.092	–
β_4	1.0	1.002	1.002	1.002	0.981	1.003	–	1.002	1.002	1.002	1.023	1.002	–	1.004	1.004	1.004	0.999	1.004	–
SD		0.163	0.163	0.163	0.238	0.163	–	0.114	0.114	0.114	0.177	0.114	–	0.096	0.096	0.096	0.145	0.095	–
SD _m		0.168	0.168	0.168	0.242	0.168	–	0.112	0.112	0.112	0.172	0.112	–	0.093	0.093	0.093	0.137	0.093	–
σ^2	9.0	9.322	9.328	9.330	29.520	9.066	–	9.415	9.418	9.419	30.491	9.275	–	9.366	9.363	9.363	29.027	9.262	–
SD		2.418	2.417	2.418	2.115	2.374	–	1.757	1.758	1.756	1.552	1.728	–	1.409	1.408	1.408	1.228	1.393	–
SD _m		2.397	2.398	2.398	–	2.343	–	1.725	1.725	1.725	–	1.700	–	1.431	1.431	1.431	–	1.416	–
c	0.2	0.238	0.238	0.238	–	0.240	–	0.202	0.202	0.202	–	0.201	–	0.193	0.193	0.193	–	0.192	–
SD		0.151	0.151	0.151	–	0.155	–	0.077	0.077	0.077	–	0.077	–	0.062	0.062	0.062	–	0.062	–
SD _m		0.134	0.134	0.134	–	0.138	–	0.075	0.075	0.075	–	0.076	–	0.061	0.061	0.061	–	0.061	–
a	1.0	0.982	0.981	0.981	–	0.996	–	0.990	0.989	0.989	–	1.000	–	1.002	1.002	1.002	–	1.009	–
SD		0.115	0.115	0.115	–	0.116	–	0.067	0.067	0.067	–	0.068	–	0.055	0.055	0.055	–	0.056	–
SD _m		0.098	0.098	0.098	–	0.101	–	0.066	0.066	0.066	–	0.067	–	0.056	0.056	0.056	–	0.056	–
b	1.0	0.977	0.977	0.977	–	1.001	–	0.986	0.986	0.985	–	1.005	–	0.993	0.993	0.993	–	1.008	–
SD		0.135	0.135	0.135	–	0.135	–	0.087	0.087	0.087	–	0.087	–	0.067	0.067	0.067	–	0.067	–
SD _m		0.127	0.127	0.127	–	0.130	–	0.083	0.083	0.083	–	0.084	–	0.067	0.067	0.067	–	0.068	–
d	0.5	0.509	0.509	0.508	–	0.513	–	0.498	0.498	0.498	–	0.500	–	0.490	0.490	0.490	–	0.491	–
SD		0.222	0.222	0.222	–	0.221	–	0.158	0.158	0.158	–	0.158	–	0.117	0.117	0.117	–	0.116	–
SD _m		0.222	0.222	0.222	–	0.224	–	0.150	0.150	0.150	–	0.150	–	0.120	0.120	0.120	–	0.120	–
e	0.5	0.502	0.501	0.501	–	0.515	–	0.497	0.497	0.497	–	0.500	–	0.483	0.483	0.483	–	0.486	–
SD		0.223	0.223	0.223	–	0.228	–	0.140	0.141	0.140	–	0.142	–	0.119	0.119	0.119	–	0.120	–
SD _m		0.217	0.217	0.217	–	0.220	–	0.145	0.145	0.145	–	0.146	–	0.118	0.118	0.118	–	0.118	–
f	0.5	0.514	0.514	0.514	–	0.524	–	0.487	0.487	0.487	–	0.491	–	0.493	0.493	0.493	–	0.496	–
SD		0.193	0.193	0.193	–	0.198	–	0.151	0.151	0.151	–	0.153	–	0.116	0.116	0.116	–	0.117	–
SD _m		0.197	0.197	0.197	–	0.200	–	0.151	0.151	0.151	–	0.151	–	0.115	0.115	0.115	–	0.115	–
ASE	–	0.825	0.828	0.829	0.823	–	–	0.577	0.577	0.578	0.576	–	–	0.450	0.447	0.447	0.450	–	–
MSPE	–	13.107	13.103	13.103	23.673	13.062	–	20.457	20.455	20.455	28.599	20.412	–	22.368	22.369	22.369	34.789	22.336	–

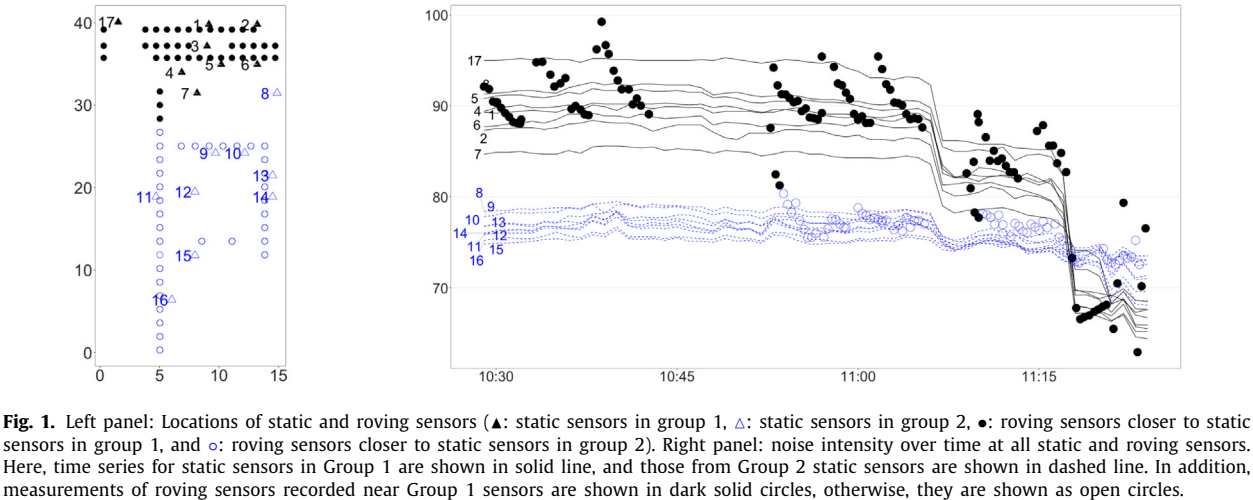


Fig. 1. Left panel: Locations of static and roving sensors (▲: static sensors in group 1, △: static sensors in group 2, ●: roving sensors closer to static sensors in group 1, and ○: roving sensors closer to static sensors in group 2). Right panel: noise intensity over time at all static and roving sensors. Here, time series for static sensors in Group 1 are shown in solid line, and those from Group 2 static sensors are shown in dashed line. In addition, measurements of roving sensors recorded near Group 1 sensors are shown in dark solid circles, otherwise, they are shown as open circles.

We apply our method to analyze this data set and summarize the parameter estimates of β_1 , β_2 and θ in Table 4, whereas the estimated temporal function $\hat{f}(t)$ and the pointwise 95% confidence intervals are plotted in Fig. 2. We

Table 4
Selected bandwidths using bimodal kernel and corresponding parameter estimates for four covariance structures: $D_1(\mathbf{s}, t) = 1$, $D_2(\mathbf{s}, t) = dt + 1$ and $D_3(\mathbf{s}, t) = dt + e(t - \kappa)_+ + 1$. Standard errors are computed based on information matrices from Theorem 1 and given in parentheses.

	$D_1(\mathbf{s}, t)$	$D_2(\mathbf{s}, t)$	$D_3(\mathbf{s}, t)$	$D_3(\mathbf{s}, t)$ (penalized)
h	0.0193	0.0193	0.0193	0.0193
Regression parameters				
β_1	-0.3922 (0.0820)	-0.4492 (0.0652)	-0.4872 (0.0608)	-0.4600 (0.0636)
β_2	0.3015 (0.0565)	0.4048 (0.0440)	0.4142 (0.0410)	0.3954 (0.0425)
Covariance parameters				
σ^2	50.8840 (8.7442)	8.8096 (1.7254)	19.0058 (3.6086)	14.7628 (2.7797)
c	0.0007 (0.0001)	0.0020 (0.0005)	0.0007 (0.0002)	0.0009 (0.0002)
c_s	0.1662 (0.0040)	0.1677 (0.0037)	0.1647 (0.0035)	0.1723 (0.0038)
c_t	0.0152 (0.0025)	0.0215 (0.0033)	0.0201 (0.0031)	0.0237 (0.0036)
d	-	1.9218 (0.2647)	-0.3070 (0.1298)	0.1982 (0.1723)
e	-	-	6.5719 (0.4883)	3.0394 (0.4177)

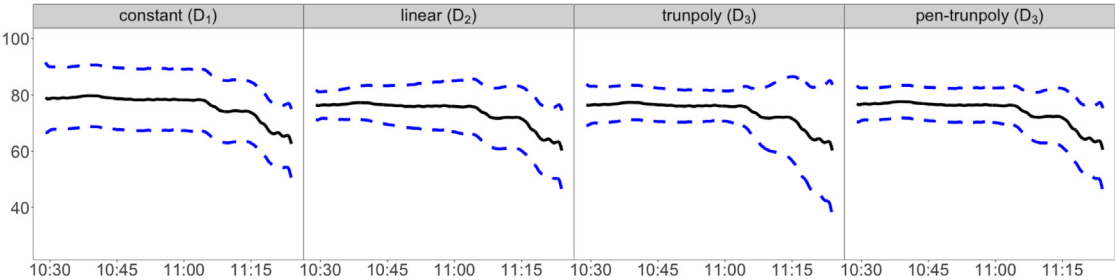


Fig. 2. Estimated temporal function $\hat{f}(t)$ (solid curve) and 95% pointwise confidence intervals (dash curves) by maximizing the profile-likelihood (6) with four covariance structures: constant $D_1(\mathbf{s}, t) = 1$; linear $D_2(\mathbf{s}, t) = dt + 1$; truncated polynomial $D_3(\mathbf{s}, t) = dt + e(t - \kappa)_+ + 1$; and maximizing a penalized profile-likelihood by adding a penalty term to (6) with D_3 .

approximate the pointwise standard deviation of $\hat{f}(t)$ by Theorem 2. The temporal function estimates $\hat{f}(t)$ under the three models D_1 , D_2 , and D_3 are quite similar; however, the pointwise confidence interval based on D_1 is much wider than those based on D_2 . For D_3 , the pointwise confidence interval is much narrower than D_1 and D_2 when t is small, however it is unusually large when t is large.

This finding is also reflected in Table 4, the estimate of the coefficient of $(t - \kappa)_+$ in D_3 (e) is unusually large. This seems like a common phenomenon in spline smoothing with truncated polynomial basis functions. To circumvent this potential issue, we consider a penalized approach [36]. That is, when maximizing the profile likelihood function (6), we consider adding an additional penalty term $-\lambda|e|$, where λ is a tuning parameter. In practice, we choose λ over a grid of λ values by minimizing the rotated residual sum of squares, for details, see Section 2.2 of the Supplementary Material. In our data analysis, $\hat{\lambda} = 20$, and the resulting parameter estimates are given in the last column of Table 4. The resulting estimate of e is much smaller, the other estimates of e are close to each other. The estimated standard deviation at each time point are plotted in Fig. C of the Supplementary Material. We notice that D_3 from the penalized approach has the smallest area under the curve. As a consequence, the 95% confidence interval of the temporal function $\hat{f}(t)$ of D_3 from the penalized approach is the narrowest compared to D_1 , D_2 and D_3 , as presented in the last panel of Fig. 2.

Finally, we consider an interpolation of the noise intensity in space and time by kriging based on D_3 with penalty. Fig. D in the Supplementary Material presents a dynamic evolution of the noise intensity maps over time and suggests a possible noise source in the upper-left corner with high noise intensity. There is also a sharp decrease of the noise intensity at 11:10:00 am when the engine was turned off even though all the sensors remain active, as well as a horizontal separation around $y = 30$ before 11:10:00 am, reflecting the wall that separates the facility [31].

CRedit authorship contribution statement

Jialuo Liu: Methodology, Software, Formal analysis, Writing - original draft. **Tingjin Chu:** Methodology, Software, Validation, Writing - review & editing. **Jun Zhu:** Conceptualization, Methodology, Writing - review & editing. **Haonan Wang:** Conceptualization, Methodology, Validation, Writing - original draft, Writing - review & editing, Project administration.

Acknowledgments

The authors would like to thank the Editor, Associate Editor for their constructive comments. The research of Haonan Wang was partially supported by National Science Foundation (NSF), United States of America grants DMS-1737795,

DMS-1923142 and CNS-1932413. This work is in part supported by the U.S. Geological Survey under a Grant/Cooperative Agreement. The views and conclusions contained in this document are those of the authors and should not be interpreted as representing the opinions or policies of the U.S. Geological Survey. Mention of trade names or commercial products does not constitute their endorsement by the U.S. Geological Survey.

Appendix A. Technical details

We use β_0 to denote the vector of true regression coefficients and θ_0 to denote the vector of true covariance parameters. We denote the log-likelihood of (β, θ) in (1), when $f(t)$ is known, as

$$\ell_0(\beta, \theta) = -(N_n/2) \log(2\pi) - (1/2) \log\{\det \Gamma(\theta)\} - (1/2)(\mathbf{y} - \mathbf{X}\beta - \mathbf{f})^\top \Gamma(\theta)^{-1}(\mathbf{y} - \mathbf{X}\beta - \mathbf{f}).$$

Let $\ell'_0(\beta) = \partial \ell_0(\beta, \theta)/\partial \beta$ and $\ell'_0(\theta) = \partial \ell_0(\beta, \theta)/\partial \theta$ denote the first-order partial derivatives of $\ell_0(\beta, \theta)$ with respect to β and θ , respectively. For ease of notation, we suppress θ in matrices relying on θ . For example, we write $\Gamma = \Gamma(\theta)$. Then, we have $\ell'_0(\beta) = \mathbf{X}^\top \Gamma^{-1}(\mathbf{y} - \mathbf{X}\beta - \mathbf{f})$ and the k th element of $\ell'_0(\theta)$ is $-(1/2)\text{tr}(\Gamma^{-1}\Gamma_k) - (1/2)(\mathbf{y} - \mathbf{X}\beta - \mathbf{f})^\top \Gamma^k(\mathbf{y} - \mathbf{X}\beta - \mathbf{f})$, where $\Gamma_k = \partial \Gamma/\partial \theta_k$ and $\Gamma^k = \partial \Gamma^{-1}/\partial \theta_k = -\Gamma^{-1}\Gamma_k\Gamma^{-1}$ for $k \in \{1, \dots, q\}$.

Further, let $\ell''_0(\beta, \beta) = \partial^2 \ell_0(\beta, \theta)/\partial \beta^2$, $\ell''_0(\theta, \theta) = \partial^2 \ell_0(\beta, \theta)/\partial \theta^2$ and $\ell''_0(\beta, \theta) = \partial^2 \ell_0(\beta, \theta)/\partial \beta \partial \theta$ denote the second-order partial derivatives with respect to β and θ . Let $\mathcal{J}_n(\beta) = E\{-\ell''_0(\beta, \beta)\}$ and $\mathcal{J}_n(\theta) = E\{-\ell''_0(\theta, \theta)\}$ denote the information matrices of β and θ , respectively. In particular, $\ell''_0(\beta, \beta) = -\mathbf{X}^\top \Gamma^{-1} \mathbf{X}$, the k th column of $\ell''_0(\beta, \theta)$ is $\mathbf{X}^\top \Gamma^k(\mathbf{y} - \mathbf{X}\beta - \mathbf{f})$, and the (k, k') th entry of $\ell''_0(\theta, \theta)$ is $-(1/2)\{\text{tr}(\Gamma^{-1}\Gamma_{kk'} + \Gamma^k\Gamma_{k'}) + (\mathbf{y} - \mathbf{X}\beta - \mathbf{f})^\top \Gamma^{kk'}(\mathbf{y} - \mathbf{X}\beta - \mathbf{f})\}$, where $\Gamma_{kk'} = \partial^2 \Gamma/\partial \theta_k \partial \theta_{k'}$ and $\Gamma^{kk'} = \partial^2 \Gamma^{-1}/\partial \theta_k \partial \theta_{k'} = \Gamma^{-1}(\Gamma_k\Gamma^{-1}\Gamma_{k'} + \Gamma_{k'}\Gamma^{-1}\Gamma_k - \Gamma_{kk'})\Gamma^{-1}$ for $k, k' \in \{1, \dots, q\}$. It can be shown that $E\{\ell''_0(\beta, \theta)\} = \mathbf{0}$, so the information matrix of η is $\mathcal{J}_n(\eta) = \text{diag}\{\mathcal{J}_n(\beta), \mathcal{J}_n(\theta)\}$, where

$$\mathcal{J}_n(\beta) = E\{-\ell''_0(\beta, \beta)\} = \mathbf{X}^\top \Gamma^{-1} \mathbf{X} \quad (13)$$

and the (k, k') th entry of $\mathcal{J}_n(\theta) = E\{-\ell''_0(\theta, \theta)\}$ is $t_{kk'}/2$ with $t_{kk'} = \text{tr}(\Gamma^{-1}\Gamma_k\Gamma^{-1}\Gamma_{k'}) = \text{tr}(\Gamma\Gamma^k\Gamma\Gamma^{k'})$.

For a matrix $\mathbf{A} = [a_{i'j'}]_{i', j'=1}^{N_n}$, we let $\mu_i(\mathbf{A})$ denote its i th largest eigenvalue, let $\|\mathbf{A}\|_2 = \mu_1(\mathbf{A})$ denote its spectral norm, let $\|\mathbf{A}\|_F = \left(\sum_{i=1}^{N_n} \sum_{j=1}^{N_n} a_{ij}^2\right)^{1/2}$ denote its Frobenius norm, let $\|\mathbf{A}\|_{\max} = \max_{i,j} |a_{ij}|$ denote its max norm, and let $\|\mathbf{A}\|_\infty = \max_{1 \leq i \leq N_n} \sum_{j=1}^{N_n} |a_{ij}|$ denote the maximum absolute column sum of the matrix. Finally, let \xrightarrow{P} denote convergence in probability and \xrightarrow{D} denote convergence in distribution, as $n \rightarrow \infty$.

The theoretical properties of the methods developed in Section 2 are established under the following additional regularity conditions.

- (C.1) There exists a nondecreasing function $Q(t)$ with $Q(0) = 0$ and $Q(1) = 1$ such that (i) $\sup_{t \in [0, 1]} |Q_{N_n}(t) - Q(t)| = \mathcal{O}(\zeta_n)$, where $Q_{N_n}(t) = N_n^{-1} \sum_{i=1}^{N_n} I(t_i \leq t)$; (ii) its first-order derivative function $q(t)$ is bounded away from zero and infinity and has continuous second partial derivatives.
- (C.2) For $j \in \{1, \dots, p\}$, there exists a function $g_j(\cdot)$ on \mathcal{T} with a bounded second derivative satisfying

$$x_j(\mathbf{s}_i, t_i) = g_j(t_i) + \phi_{ij}, \quad i \in \{1, \dots, N_n\},$$

where $\{\phi_{ij}\}$ is a sequence of real numbers such that

$$\lim_{n \rightarrow \infty} N_n^{-1} \Phi^\top \Gamma^{-1} \Phi = \Pi,$$

where $\Phi_i = (\phi_{i1}, \dots, \phi_{iN_n})^\top$, $\Phi = (\Phi_1, \dots, \Phi_p)$, and Π is a positive definite matrix. In addition, for $j \in \{1, \dots, p\}$, $\limsup_{n \rightarrow \infty} (1/a_n) \max_{1 \leq k \leq N_n} \left| \sum_{m=1}^k \phi_{imj} \right| < \infty$ for all permutations (i_1, \dots, i_{N_n}) of $(1, \dots, N_n)$, where $a_n = N_n^{1/2} \log N_n$.

- (C.3) The temporal function $f(t)$ is twice differentiable with a bounded second-order derivative on \mathcal{T} .
- (C.4) The kernel $K(\cdot)$ is a symmetric, nonnegative, and bounded function with a compact support in \mathbb{R} and with a bounded first-order derivative.
- (C.5) The bandwidth h satisfies $h \rightarrow 0$, $N_n h^4 \rightarrow \infty$, $N_n h^8 \rightarrow 0$ and $\zeta_n h^{-1} \rightarrow 0$ as $n \rightarrow \infty$.
- (C.6) For $k \in \{1, \dots, q\}$, $\|\Gamma_k\|_F^{-2} \leq D_k N_n^{-1/2-\iota}$ for some $\iota > 0$ and $D_k > 0$.
- (C.7) It holds that $\|\Gamma^{-1}\|_2 < C^* < \infty$ for some constant C^* .
- (C.8) $\lim_{n \rightarrow \infty} N_n^{-1} \mathcal{J}_n(\theta) \rightarrow \mathcal{I}_0(\theta)$, where $\mathcal{I}_0(\theta)$ is non-singular.
- (C.9) Given $t \in (0, 1)$, there exists a 2×2 matrix Δ_t , such that $(N_n^{-1} h) \mathbf{k}_t^\top \Gamma \mathbf{k}_t \rightarrow q(t)^2 \Delta_t$, where $\mathbf{k}_t = \{K_h(t_i - t)\}_{i,j=1}^{N_n, 2}$ is an $N_n \times 2$ matrix.
- (C.10) Define $g(\mathbf{s}, t) = g(\mathbf{0}, \mathbf{0}, \mathbf{s}, t)$. Assume $g(\mathbf{s}, t)$ satisfies $|g(\mathbf{s}, t) - g(\mathbf{s}', t')| \leq C_1 \|\mathbf{s} - \mathbf{s}'\| + C_2 |t - t'|$ for all $(\mathbf{s}, t), (\mathbf{s}', t') \in \mathcal{R} \times \mathcal{T}$, where C_1, C_2 are positive constants.
- (C.11) There exist two positive nonincreasing functions γ_0 and γ_1 such that $|\gamma_n((\mathbf{s}, t), (\mathbf{s} + \mathbf{u}_1/A_n, t + u_2/B_n))| \leq \gamma_0(\|\mathbf{u}_1\|) \gamma_1(|u_2|)$ for all n and $\|\mathbf{u}_1\|, |u_2| \in [0, \infty)$ such that $(\mathbf{s}, t), (\mathbf{s} + \mathbf{u}_1/A_n, t + u_2/B_n) \in \mathcal{R} \times \mathcal{T}$. In addition, $\int_0^\infty u^{d-1} \gamma_0(u) du < \infty$ and $\int_0^\infty \gamma_1(u) du < \infty$.

(C.12) The covariance function $\gamma_n(\cdot, \cdot; \theta)$ is bounded and is twice continuously differentiable with respect to θ in an open set.

(C.13) There exist two positive nonincreasing functions γ_2 and γ_3 such that

$$\max\{|\gamma_{n,k}((\mathbf{s}, t), (\mathbf{s} + \mathbf{u}_1/A_n, t + u_2/B_n))|, |\gamma_{n,kk'}((\mathbf{s}, t), (\mathbf{s} + \mathbf{u}_1/A_n, t + u_2/B_n))|\} \leq \gamma_2(\|\mathbf{u}_1\|)\gamma_3(|u_2|)$$

for all n and $\|\mathbf{u}_1\|, |u_2| \in [0, \infty)$ with $(\mathbf{s}, t), (\mathbf{s} + \mathbf{u}_1/A_n, t + u_2/B_n) \in \mathcal{R} \times \mathcal{T}$ and $1 \leq k, k' \leq q$. Further, $\int_0^\infty u^{d-1} \gamma_2(u) du < \infty$ and $\int_0^\infty \gamma_3(u) du < \infty$.

In the following proofs, we suppress n in ${}^n t_{kk'}, {}^n a_{kk'}, {}^n \Gamma, {}^n \Gamma_k, {}^n \Gamma_{kk'}, \mathbf{I}_n, \mathbf{A}_n, {}^n \hat{\boldsymbol{\eta}}, {}^n \hat{\boldsymbol{\beta}}$ and ${}^n \hat{\boldsymbol{\theta}}$ for ease of notation.

Remarks. (C.1) is a condition on fixed time points for the spatio-temporal sampling design. (C.2) is a mild assumption about the relationship between the fixed design points and $\{t_i\}$ in the partially linear model, which is similar to Assumption 2.2 (i) in Gao and Liang [17]. (C.3)–(C.5) are common assumptions in kernel smoothing. (C.3) ensures the smoothness of the temporal function [27,43]. (C.4) is a standard assumption for kernel functions and can be relaxed further such that $K(t)$ satisfies a Lipschitz condition $|K(t) - K(t')| \leq c|t - t'|$ for any $t, t' \in \mathbb{R}$ and some $c > 0$. In addition, (C.5) is a condition for the rate of bandwidth with respect to N_n and ζ_n . (C.6) assures that the first-order partial derivatives of the covariance matrix have a higher order than root- N_n . (C.7) imposes a lower bound on the smallest eigenvalue of the covariance matrix. (C.8) guarantees that the growth of the information matrix is at the rate of the total sample size [6]. Moreover, (C.9) is an assumption for the fixed sampling design under the spatio-temporal dependence. Finally, (C.10)–(C.13) are regularity conditions for locally stationary processes. In (C.11), the covariance function of locally stationary processes is bounded by a product of two functions, whose integrals are finite.

A Remark on Assumption (C.2) Here, we will show that if $\|{}^n \Gamma^{-1}\|_\infty = \mathcal{O}(1)$, we have $\mathbf{X}^\top \Gamma^{-1} \mathbf{X} \succeq \Phi^\top \Gamma^{-1} \Phi$, where $\mathbf{A} \succeq \mathbf{B}$ if $\mathbf{A} - \mathbf{B}$ is positive semi-definite. To see this, we write

$$\mathbf{X}^\top \Gamma^{-1} \mathbf{X} = \mathbf{G}^\top \Gamma^{-1} \mathbf{G} + \mathbf{G}^\top \Gamma^{-1} \Phi + \Phi^\top \Gamma^{-1} \mathbf{G} + \Phi^\top \Gamma^{-1} \Phi.$$

Since $\|{}^n \Gamma^{-1}\|_\infty = \mathcal{O}(1)$, $\mathbf{G}^\top \Gamma^{-1}$ is uniformly bounded elementwise. Together with (C.2), $\mathbf{G}^\top \Gamma^{-1} \Phi = \mathcal{O}(N_n^{1/2} \log N_n)$. Recall that $\lim_{n \rightarrow \infty} N_n^{-1} \Phi^\top \Gamma^{-1} \Phi = \Pi$. Thus, $\mathbf{G}^\top \Gamma^{-1} \Phi$ is dominated by $\Phi^\top \Gamma^{-1} \Phi$, and thus,

$$\mathbf{X}^\top \Gamma^{-1} \mathbf{X} \succeq \Phi^\top \Gamma^{-1} \Phi,$$

in which the equality holds if $g(\cdot) = 0$. This result indicates that the asymptotic variances of $\hat{\boldsymbol{\beta}}$ in the partially linear model are greater than those in the simple linear regression model.

In the following Lemmas 1–6, we generalize some classical results for random sampling designs [14] to fixed sampling design, which will be used in the proofs of Theorems 1–3.

Lemma 1. Under Assumptions (C.1), (C.4) and (C.5), for $k \geq 0$,

$$\sup_{t \in [0,1]} |v_{k,t} - N_n \mu_{k,t} q(t)| = \mathcal{O}(N_n h + N_n \zeta_n h^{-1}),$$

where $v_{k,t} = h^{-k} \sum_{i=1}^{N_n} (t_i - t)^k K_h(t_i - t)$, $K_h(t) = (1/h)K(t/h)$,

$$\mu_{k,t} = \begin{cases} \int_{-t/h}^\infty x^k K(x) dx, & \text{if } t < Mh, \\ \int_{-\infty}^\infty x^k K(x) dx := \mu_k, & \text{if } Mh \leq t \leq 1 - Mh, \\ \int_{-\infty}^{(1-t)/h} x^k K(x) dx, & \text{if } t > 1 - Mh, \end{cases}$$

and $[-M, M]$ is the compact support of $K(\cdot)$.

Proof. For any $t \in [0, 1]$,

$$\begin{aligned} |v_{k,t} - N_n \mu_{k,t} q(t)| &\leq \left| N_n h^{-k} \int_0^1 (z - t)^k K_h(z - t) d(Q_{N_n} - Q)(z) \right| \\ &\quad + \left| N_n h^{-k} \int_0^1 (z - t)^k K_h(z - t) dQ(z) - N_n \mu_{k,t} q(t) \right| \equiv (I_{1,1}) + (I_{1,2}), \end{aligned}$$

where

$$\begin{aligned} (I_{1,1}) &= N_n h^{-k} \left| \int_0^1 (z - t)^k K_h(z - t) d(Q_{N_n} - Q)(z) \right| \\ &= N_n h^{-k} \left| (z - t)^k K_h(z - t) (Q_{N_n} - Q)(z) \Big|_0^1 - \int_0^1 (Q_{N_n} - Q)(z) [(z - t)^k K_h(z - t)]' dz \right| \end{aligned}$$

$$\begin{aligned}
&= N_n h^{-2} \left| \int_0^1 (Q_{N_n} - Q)(z) k \left(\frac{z-t}{h} \right)^{k-1} K \left(\frac{z-t}{h} \right) dz + \int_0^1 (Q_{N_n} - Q)(z) \left(\frac{z-t}{h} \right)^k K' \left(\frac{z-t}{h} \right) dz \right| \\
&\leq N_n h^{-1} \sup_{z \in [0,1]} |(Q_{N_n} - Q)(z)| \left(\int_{-M}^M |k u^{k-1} K(u)| du + \int_{-M}^M |u^k K'(u)| du \right).
\end{aligned}$$

The second equality uses the fact that $(Q_{N_n} - Q)(1) = (Q_{N_n} - Q)(0) = 0$. By (C.1), $\int_{-M}^M |k u^{k-1} K(u)| du + \int_{-M}^M |u^k K'(u)| du = \mathcal{O}(1)$. Together with (C.4), $(I_{1,1}) = \mathcal{O}(N_n \zeta_n h^{-1})$. For $(I_{1,2})$,

$$\begin{aligned}
(I_{1,2}) &= \sup_{t \in [0,1]} \left| N_n h^{-k} \int_0^1 (z-t)^k K_h(z-t) dQ(z) - N_n q(t) \mu_{k,t} \right| \\
&= \sup_{t \in [0,1]} \left| N_n \int_{-\frac{t}{h}}^{\frac{1-t}{h}} u^k K(u) \left(q(t) + q'(t)uh + \frac{q''(\tilde{t})u^2 h^2}{2} \right) du - N_n q(t) \mu_{k,t} \right| \\
&\leq N_n h \left\{ \sup_{t \in [0,1]} |q'(t)| \int |u^{k+1} K(u)| du + (h/2) \sup_{t \in [0,1]} |q''(t)| \int |u^{k+2} K(u)| du \right\} = \mathcal{O}(N_n h),
\end{aligned}$$

where $\tilde{t} \in [t, t + uh]$. Thus, Lemma 1 holds. \square

Lemma 2. Under Assumptions (C.1) and (C.3)–(C.5), $\sup_{t \in [0,1]} |\omega_1(t)\mathbf{f} - f(t)| = \mathcal{O}(h^2)$.

Proof. First, straightforward calculation yields $\mathbf{D}_t^\top \mathbf{K}_t \mathbf{D}_t = \begin{pmatrix} v_{0,t} & v_{1,t} \\ v_{1,t} & v_{2,t} \end{pmatrix}$. By Lemma 1, uniformly on $[0, 1]$, we have $v_{0,t} = N_n q(t) \mu_{0,t} + \mathcal{O}(N_n h + N_n \zeta_n h^{-1})$, $v_{1,t} = N_n q(t) \mu_{1,t} + \mathcal{O}(N_n h + N_n \zeta_n h^{-1})$ and $v_{2,t} = N_n q(t) \mu_{2,t} + \mathcal{O}(N_n h + N_n \zeta_n h^{-1})$. In addition, notice that

$$(1, 0)(\mathbf{D}_t^\top \mathbf{K}_t \mathbf{D}_t)^{-1} = \begin{pmatrix} \frac{v_{2,t}}{v_{0,t}v_{2,t} - v_{1,t}^2}, \frac{-v_{1,t}}{v_{0,t}v_{2,t} - v_{1,t}^2} \end{pmatrix}.$$

Thus,

$$\begin{aligned}
\frac{v_{2,t}}{v_{0,t}v_{2,t} - v_{1,t}^2} &= N_n^{-1}(q(t))^{-1} \frac{\mu_{2,t}}{\mu_{0,t}\mu_{2,t} - \mu_{1,t}^2} + \mathcal{O}(N_n^{-1}h + N_n^{-1}\zeta_n h^{-1}), \\
\frac{-v_{1,t}}{v_{0,t}v_{2,t} - v_{1,t}^2} &= N_n^{-1}(q(t))^{-1} \frac{\mu_{1,t}}{\mu_{0,t}\mu_{2,t} - \mu_{1,t}^2} + \mathcal{O}(N_n^{-1}h + N_n^{-1}\zeta_n h^{-1})
\end{aligned}$$

uniformly on $[0, 1]$.

Recall that $\omega_1(t)\mathbf{f} - f(t) = (1, 0)(\mathbf{D}_t^\top \mathbf{K}_t \mathbf{D}_t)^{-1} \mathbf{D}_t^\top \mathbf{K}_t \mathbf{f} - f(t)$. A Taylor's expansion yields $f(t_i) = f(t) + f'(t)(t_i - t) + 1/2 f''(\xi_i)(t_i - t)^2$, where ξ_i is between t and t_i . Thus,

$$\begin{aligned}
\omega_1(t)\mathbf{f} - f(t) &= (1, 0)(\mathbf{D}_t^\top \mathbf{K}_t \mathbf{D}_t)^{-1} \mathbf{D}_t^\top \mathbf{K}_t \mathbf{D}_t (f(t), hf'(t))^\top + (1/2)(1, 0)(\mathbf{D}_t^\top \mathbf{K}_t \mathbf{D}_t)^{-1} \mathbf{D}_t^\top \mathbf{K}_t \mathbf{d}_\xi - f(t) \\
&= (1/2)(1, 0)(\mathbf{D}_t^\top \mathbf{K}_t \mathbf{D}_t)^{-1} \mathbf{D}_t^\top \mathbf{K}_t \mathbf{d}_\xi,
\end{aligned}$$

where $\mathbf{d}_\xi = (f''(\xi_1)(t_1 - t)^2, \dots, f''(\xi_{N_n})(t_{N_n} - t)^2)^\top$. In addition,

$$\sup_{t \in [0,1]} |(1, 0)(\mathbf{D}_t^\top \mathbf{K}_t \mathbf{D}_t)^{-1} \mathbf{D}_t^\top \mathbf{K}_t \mathbf{d}_\xi| \leq \max_{x \in [0,1]} |f''(x)| \sup_{t \in [0,1]} \left(\left| \frac{v_{2,t}}{v_{0,t}v_{2,t} - v_{1,t}^2} \right| + \left| \frac{v_{1,t}v_{3,t}}{v_{0,t}v_{2,t} - v_{1,t}^2} \right| \right) h^2 = \mathcal{O}(h^2).$$

Thus, Lemma 2 holds. \square

Lemma 3. Suppose that Assumptions (C.1) and (C.3)–(C.5) hold. For any random vector $\boldsymbol{\varepsilon}$ of zero mean,

$$\sup_{t \in [0,1]} |\omega_1(t)\boldsymbol{\varepsilon}| = \mathcal{O}_p \left\{ \left(\frac{\log N_n}{N_n h} \right)^{1/2} \right\}.$$

Proof. For a random vector $\boldsymbol{\varepsilon} = (\varepsilon_1, \dots, \varepsilon_{N_n})^\top$, we have $\omega_1(t)\boldsymbol{\varepsilon} = (1, 0)(\mathbf{D}_t^\top \mathbf{K}_t \mathbf{D}_t)^{-1} \mathbf{D}_t^\top \mathbf{K}_t \boldsymbol{\varepsilon} = (I_{3,1}) - (I_{3,2})$, where $(I_{3,1}) = \frac{v_{2,t}}{v_{0,t}v_{2,t} - v_{1,t}^2} \sum_{i=1}^{N_n} K_h(t_i - t) \varepsilon_i$ and $(I_{3,2}) = \frac{v_{1,t}}{v_{0,t}v_{2,t} - v_{1,t}^2} \sum_{i=1}^{N_n} K_h(t_i - t)(t_i - t) h^{-1} \varepsilon_i$.

Note that $(I_{3,1}) = \frac{v_{0,t} v_{2,t}}{v_{0,t} v_{2,t} - v_{1,t}^2} \frac{\sum_{i=1}^{N_n} K_h(t_i - t) \varepsilon_i}{\sum_{i=1}^{N_n} K_h(t_i - t)}$, by Lemma 5,

$$\sup_{t \in [0,1]} \left| \frac{\sum_{i=1}^{N_n} K_h(t_i - t) \varepsilon_i}{\sum_{i=1}^{N_n} K_h(t_i - t)} \right| = \mathcal{O}_p \left\{ \left(\frac{\log N_n}{N_n h} \right)^{1/2} \right\}, \quad \sup_{t \in [0,1]} \left| \frac{\sum_{i=1}^{N_n} (t_i - t) h^{-1} K_h(t_i - t) \varepsilon_i}{\sum_{i=1}^{N_n} K_h(t_i - t)} \right| = \mathcal{O}_p \left\{ \left(\frac{\log N_n}{N_n h} \right)^{1/2} \right\}.$$

Following similar arguments in Lemma 2, $\sup_{t \in [0,1]} |(I_{3,1})| = \mathcal{O}_p \left\{ \left(\frac{\log N_n}{N_n h} \right)^{1/2} \right\}$ and $\sup_{t \in [0,1]} |(I_{3,2})| = \mathcal{O}_p \left\{ \left(\frac{\log N_n}{N_n h} \right)^{1/2} \right\}$. \square

Lemma 4. Under Assumptions (C.1) and (C.3)–(C.5),

$$\sup_{t \in [0,1]} |\tilde{f}(t) - f(t)| = \mathcal{O}_p \left\{ h^2 + \left(\frac{\log N_n}{N_n h} \right)^{1/2} \right\},$$

where $\tilde{f}(t) = \omega_1(t) \mathbf{y}^* = \omega_1(t) (\mathbf{y} - \mathbf{X} \beta)$.

Proof. First, $\tilde{f}(t) - f(t) = \omega_1(t) \{ \mathbf{f} + \boldsymbol{\varepsilon} - f(t) \mathbf{1}_{N_n} \}$ since $\omega_1(t) \mathbf{1}_{N_n} - 1 = 0$. Next, we have $\sup_{t \in [0,1]} |\tilde{f}(t) - f(t)| \leq \sup_{t \in [0,1]} |\omega_1(t) (\mathbf{f} - f(t) \mathbf{1}_{N_n})| + \sup_{t \in [0,1]} |\omega_1(t) \boldsymbol{\varepsilon}|$. The desired result follows from Lemmas 2 and 3. \square

Lemma 5. Suppose Assumptions (C.1), (C.4) and (C.5) hold. For any random vector $\boldsymbol{\varepsilon}$ of zero mean,

$$\sup_{t \in [0,1]} \left| \frac{\sum_{i=1}^{N_n} (t_i - t)^j h^{-j} K_h(t_i - t) \varepsilon_i}{\sum_{i=1}^{N_n} K_h(t_i - t)} \right| = \mathcal{O}_p \left\{ \left(\frac{\log N_n}{N_n h} \right)^{1/2} \right\}, \quad j \in \{0, 1\}.$$

Proof. Let I_k be the interval centered at c_k with the length $\iota_{N_n} = \{\log N_n / (N_n h)\}^{1/2} h^{3+j}$. There exist $r_{N_n} = \lfloor \iota_{N_n}^{-1} \rfloor + 1$ intervals satisfying $[0, 1] \subset \bigcup_{k=1}^{r_{N_n}} I_k$. First, $\sup_{t \in [0,1]} \left| \frac{\sum_{i=1}^{N_n} (t_i - t)^j h^{-j} K_h(t_i - t) \varepsilon_i}{\sum_{i=1}^{N_n} K_h(t_i - t)} \right| \leq (I_{6,1}) + (I_{6,2})$, where

$$(I_{6,1}) = \max_{1 \leq k \leq r_{N_n}} \sup_{t \in I_k} \left| \frac{\sum_{i=1}^{N_n} (t_i - t)^j K_h(t_i - t) \varepsilon_i}{h^j \sum_{i=1}^{N_n} K_h(t_i - t)} - \frac{\sum_{i=1}^{N_n} (t_i - c_k)^j K_h(t_i - c_k) \varepsilon_i}{h^j \sum_{i=1}^{N_n} K_h(t_i - c_k)} \right|,$$

$$(I_{6,2}) = \max_{1 \leq k \leq r_{N_n}} \left| \frac{\sum_{i=1}^{N_n} (t_i - c_k)^j h^{-j} K_h(t_i - c_k) \varepsilon_i}{\sum_{i=1}^{N_n} K_h(t_i - c_k)} \right|.$$

For $(I_{6,1})$,

$$(I_{6,1}) \leq \max_{1 \leq k \leq r_{N_n}} \sup_{t \in I_k} \left| \frac{1}{v_{0,t}} \left[\sum_{i=1}^{N_n} \bar{K}(t, t_i, c_k) (t_i - t)^j h^{-j} \varepsilon_i \right] \right|$$

$$+ \max_{1 \leq k \leq r_{N_n}} \sup_{t \in I_k} \left| \frac{1}{v_{0,t}} \left[\sum_{i=1}^{N_n} \{ (t_i - t)^j - (t_i - c_k)^j \} h^{-j} K_h(t_i - c_k) \varepsilon_i \right] \right|$$

$$+ \max_{1 \leq k \leq r_{N_n}} \sup_{t \in I_k} \left| \frac{1}{v_{0,t} v_{0,c_k}} \sum_{i=1}^{N_n} \bar{K}(t, t_i, c_k) \sum_{i=1}^{N_n} (t_i - c_k)^j h^{-j} K_h(t_i - c_k) \varepsilon_i \right|$$

$$= (I_{6,1A}) + (I_{6,1B}) + (I_{6,1C}),$$

where $\bar{K}(t, t_i, c_k) = K_h(t_i - t) - K_h(t_i - c_k)$.

By Lemma 1, it can be shown that $\sup_{t \in [0,1]} |v_{0,t}^{-1}| = \mathcal{O}(N_n^{-1})$. In addition, by (C.4), for any $t \in I_k$, $|\bar{K}(t, t_i, c_k)| \leq h^{-1} \max_{x \in \mathbb{R}} |K'(x)| \left| \frac{t_i - t}{h} - \frac{t_i - c_k}{h} \right| = \mathcal{O}(h^{-2} \iota_{N_n})$. Therefore,

$$(I_{6,1A}) = \max_{1 \leq k \leq r_{N_n}} \sup_{t \in I_k} \left| \frac{1}{v_{0,t}} \left[\sum_{i=1}^{N_n} \bar{K}(t, t_i, c_k) (t_i - t)^j h^{-j} \varepsilon_i \right] \right|$$

$$\leq \max_{1 \leq k \leq r_{N_n}} \sup_{t \in I_k} |\bar{K}(t, t_i, c_k)| \max_{1 \leq k \leq r_{N_n}} \sup_{t \in I_k} \left\{ \frac{1}{v_{0,t}} \sum_{i=1}^{N_n} |(t_i - t)^j h^{-j} \varepsilon_i| \right\} = \mathcal{O}(N_n^{-1} \iota_{N_n} h^{-2-j}) \sum_{i=1}^{N_n} |\varepsilon_i|.$$

We further note that

$$(I_{6,1B}) = \max_{1 \leq k \leq r_{N_n}} \sup_{t \in I_k} \left| \frac{1}{v_{0,t}} \left[\sum_{i=1}^{N_n} \{(t_i - t)^j - (t_i - c_k)^j\} h^{-j} K_h(t_i - c_k) \varepsilon_i \right] \right|$$

$$= \begin{cases} 0, & \text{if } j = 0, \\ \mathcal{O}(N_n^{-1} \iota_{N_n} h^{-1-j}) \sum_{i=1}^{N_n} |\varepsilon_i|, & \text{if } j = 1. \end{cases}$$

Moreover,

$$(I_{6,1C}) = \max_{1 \leq k \leq r_{N_n}} \sup_{t \in I_k} \left| \frac{1}{v_{0,t} v_{0,c_k}} \sum_{i=1}^{N_n} \bar{K}(t, t_i, c_k) \sum_{i=1}^{N_n} (t_i - c_k)^j h^{-j} K_h(t_i - c_k) \varepsilon_i \right|$$

$$\leq \max_{1 \leq k \leq r_{N_n}} \sup_{t \in I_k} \left\{ \frac{1}{v_{0,t} v_{0,c_k}} \sum_{i=1}^{N_n} |\bar{K}(t, t_i, c_k)| \right\} \max_{1 \leq k \leq r_{N_n}} \left\{ \sum_{i=1}^{N_n} |(t_i - c_k)^j h^{-j} \varepsilon_i| K_h(t_i - c_k) \right\}$$

$$= \mathcal{O}(N_n^{-2}) \mathcal{O}(N_n \iota_{N_n} h^{-2}) \mathcal{O}(h^{-1-j}) \sum_{i=1}^{N_n} |\varepsilon_i| = \mathcal{O}(N_n^{-1} \iota_{N_n} h^{-3-j}) \sum_{i=1}^{N_n} |\varepsilon_i|.$$

Since $\sum_{i=1}^{N_n} |\varepsilon_i| = \mathcal{O}_p(N_n)$, $(I_{6,1}) = \mathcal{O}_p \left\{ \left(\frac{\log N_n}{N_n h} \right)^{1/2} \right\}$.

For $(I_{6,2})$, let $\mathbf{e} = \Gamma^{-1/2} \boldsymbol{\varepsilon}$ be a sequence of iid $N(0, 1)$, and we have $\sum_{i=1}^{N_n} (t_i - t)^j K_h(t_i - t) \varepsilon_i = h^j \mathbf{k}_{t,j+1}^\top \Gamma^{1/2} \mathbf{e}$, where $\mathbf{k}_{t,j+1}$ is the $(j+1)$ th column of \mathbf{k}_t , $j \in \{0, 1\}$. For any $\lambda > 0$ and $t \in [0, 1]$, by Bernstein inequality,

$$P \left(\left| \sum_{i=1}^{N_n} (t_i - t)^j K_h(t_i - t) \varepsilon_i \right| > 2\lambda v_{0,t} \left(\frac{\log N_n}{N_n h} \right)^{1/2} \right) < \exp \left\{ \frac{-\lambda^2 v_{0,t}^2 \frac{\log N_n}{N_n h}}{h^{2j} \mathbf{k}_{t,j+1}^\top \Gamma \mathbf{k}_{t,j+1}} \right\}.$$

In addition, we have $h^{2j} \mathbf{k}_{t,j+1}^\top \Gamma \mathbf{k}_{t,j+1} \leq \|\Gamma\|_2 \sum_{i=1}^{N_n} (t_i - t)^{2j} K_h(t_i - t)^2 \leq \|\Gamma\|_2 \sum_{i=1}^{N_n} K_h(t_i - t)^2$. By similar arguments as in [Lemma 1](#), we can show that, $\sup_{t \in [0,1]} \left| \sum_{i=1}^{N_n} K_h(t_i - t)^2 \right| = \mathcal{O}(N_n h^{-1})$. From [Lemma 1](#), we also have $\inf_{t \in [0,1]} v_{0,t}^2 = \mathcal{O}(N_n^2)$, and therefore, by choosing a large enough λ ,

$$\sup_{t \in [0,1]} P \left(\left| \sum_{i=1}^{N_n} (t_i - t)^j K_h(t_i - t) \varepsilon_i \right| > \lambda v_{0,t} \left(\frac{\log N_n}{N_n h} \right)^{1/2} \right) = \mathcal{O}(N_n^{-2}).$$

Since

$$P \left((I_{6,2}) > \lambda \left(\frac{\log N_n}{N_n h} \right)^{1/2} \right) \leq \sum_{k=1}^{r_{N_n}} P \left(\left| \frac{\sum_{i=1}^{N_n} (t_i - c_k)^j K_h(t_i - c_k) \varepsilon_i}{v_{0,c_k}} \right| > \lambda \left(\frac{\log N_n}{N_n h} \right)^{1/2} \right)$$

$$= \mathcal{O}(r_{N_n} N_n^{-2}) = \mathcal{O}(1),$$

we have $(I_{6,2}) = \mathcal{O}_p \left\{ \left(\frac{\log N_n}{N_n h} \right)^{1/2} \right\}$. Thus, we have the result. \square

Lemma 6. Suppose that Assumptions (C.1)–(C.5) hold, we have $\sup_{t \in [0,1]} |\boldsymbol{\omega}(t) \mathbf{X}| = \mathcal{O}(1)$.

Proof. Using similar arguments as in [Lemma 2](#), $v_{k,t}/(v_{0,t} v_{2,t} - v_{1,t}^2) = \mathcal{O}(N_n^{-1})$, for $k \in \{0, 1, 2\}$. The i th element of the first row of $\boldsymbol{\omega}(t)$ is $v_{2,t}/(v_{0,t} v_{2,t} - v_{1,t}^2) K_h(t_i - t) - v_{1,t}/(v_{0,t} v_{2,t} - v_{1,t}^2) K_h(t_i - t)(t_i - t)/h = \mathcal{O}(N_n^{-1} h^{-1})$. Similarly, the i th element of the second row of $\boldsymbol{\omega}(t)$ is $\mathcal{O}(N_n^{-1} h^{-1})$. Thus, $\boldsymbol{\omega}(t) \boldsymbol{\phi}_j = \mathcal{O}(N_n^{-1/2} h^{-1} \log N_n)$. Using similar arguments in [Lemma 2](#), we obtain $(0, 1) \boldsymbol{\omega}(t) \mathbf{g}_j - h g_j'(t) = (0, 1/2) (\mathbf{D}_t^\top \mathbf{K}_t \mathbf{D}_t)^{-1} \mathbf{D}_t^\top \mathbf{K}_t \mathbf{d}_{\xi,j} = \mathcal{O}(h^2)$, where ξ_i is between t and t_i and $\mathbf{d}_{\xi,j} = (g''(\xi_1)(t_1 - t)^2, \dots, g''(\xi_{N_n})(t_{N_n} - t)^2)^\top$. Thus, $\boldsymbol{\omega}(t) \mathbf{X}_j = \boldsymbol{\omega}(t) (\boldsymbol{\phi}_j + \mathbf{g}_j) = \mathcal{O}(1)$. \square

Proof of Theorem 1. By Mardia and Marshall [32], the convergence property of $\ell'_0(\boldsymbol{\beta})$, $\ell'_0(\boldsymbol{\theta})$, $\ell''_0(\boldsymbol{\beta}, \boldsymbol{\beta})$, $\ell''_0(\boldsymbol{\beta}, \boldsymbol{\theta})$ and $\ell''_0(\boldsymbol{\theta}, \boldsymbol{\theta})$ can be established. By (C.6)–(C.7), together with proof of Theorem 1 in Chu et al. [5], we have

$$N_n^{-1/2} \ell'_0(\boldsymbol{\theta}) \xrightarrow{D} N(\mathbf{0}, \mathcal{I}_0(\boldsymbol{\theta})), N_n^{-1} \ell''_0(\boldsymbol{\theta}, \boldsymbol{\theta}) \xrightarrow{P} -\mathcal{I}_0(\boldsymbol{\theta}).$$

Under (C.1)–(C.9), we first show the following results

$$N_n^{-1/2} \ell'(\boldsymbol{\beta}) \xrightarrow{D} N(\mathbf{0}, \boldsymbol{\Pi}), N_n^{-1/2} \ell'(\boldsymbol{\theta}) \xrightarrow{D} N(\mathbf{0}, \mathcal{I}_0(\boldsymbol{\theta})), N_n^{-1} \ell''(\boldsymbol{\beta}, \boldsymbol{\beta}) \xrightarrow{P} -\boldsymbol{\Pi}, \quad (14a)$$

$$N_n^{-1} \ell''(\boldsymbol{\theta}, \boldsymbol{\theta}) \xrightarrow{P} -\mathcal{I}_0(\boldsymbol{\theta}), N_n^{-1} \ell''(\boldsymbol{\beta}, \boldsymbol{\theta}) \xrightarrow{P} \mathbf{0}. \quad (14b)$$

First, straightforward calculation yields

$$\begin{aligned}\ell'(\boldsymbol{\beta}) &= \mathbf{X}^\top (\mathbf{I} - \mathbf{S})^\top \mathbf{\Gamma}^{-1} (\mathbf{I} - \mathbf{S}) (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}) \\ &= \mathbf{X}^\top (\mathbf{I} - \mathbf{S})^\top \mathbf{\Gamma}^{-1} (\mathbf{I} - \mathbf{S}) \boldsymbol{\varepsilon} + \mathbf{X}^\top (\mathbf{I} - \mathbf{S})^\top \mathbf{\Gamma}^{-1} (\mathbf{I} - \mathbf{S}) \mathbf{f} \equiv (I_1) + (I_2).\end{aligned}$$

By Assumption (C.2),

$$\begin{aligned}(I_1) &= N_n^{-1/2} \{\boldsymbol{\phi}_j + \mathbf{g}_j\}^\top (\mathbf{I} - \mathbf{S})^\top \mathbf{\Gamma}^{-1} (\mathbf{I} - \mathbf{S}) \boldsymbol{\varepsilon} \\ &= N_n^{-1/2} \boldsymbol{\phi}_j^\top \mathbf{\Gamma}^{-1} \boldsymbol{\varepsilon} + N_n^{-1/2} \boldsymbol{\phi}_j^\top \mathbf{\Gamma}^{-1} \mathbf{S} \boldsymbol{\varepsilon} + N_n^{-1/2} \mathbf{g}_j^\top (\mathbf{I} - \mathbf{S})^\top \mathbf{\Gamma}^{-1} \mathbf{S} \boldsymbol{\varepsilon} + \\ &\quad N_n^{-1/2} \boldsymbol{\phi}_j^\top \mathbf{S}^\top \mathbf{\Gamma}^{-1} \boldsymbol{\varepsilon} + N_n^{-1/2} \mathbf{g}_j^\top (\mathbf{I} - \mathbf{S})^\top \mathbf{\Gamma}^{-1} \boldsymbol{\varepsilon} + N_n^{-1/2} \boldsymbol{\phi}_j^\top \mathbf{S}^\top \mathbf{\Gamma}^{-1} \mathbf{S} \boldsymbol{\varepsilon} \\ &\equiv (I_{11}) + (I_{12}) + (I_{13}) + (I_{14}) + (I_{15}) + (I_{16}).\end{aligned}$$

For (I_{11}) , it can be shown that $N_n^{-1/2} \boldsymbol{\phi}_j^\top \mathbf{\Gamma}^{-1} \boldsymbol{\varepsilon} \xrightarrow{D} N(\mathbf{0}, N_n^{-1} \boldsymbol{\phi}_j^\top \mathbf{\Gamma}^{-1} \boldsymbol{\phi}_j)$. In addition, we have

$$\begin{aligned}N_n^{-1} \mathbb{E}(\boldsymbol{\phi}_j^\top \mathbf{\Gamma}^{-1} \mathbf{S} \boldsymbol{\varepsilon})^2 &\leq N_n^{-1} \|\mathbf{\Gamma}^{-1}\|_2^2 \mathbb{E}|\boldsymbol{\phi}_j^\top \mathbf{S} \boldsymbol{\varepsilon}|^2 = \mathcal{O}((\log N_n)^3 N_n^{-1} h^{-1}), \\ N_n^{-1} \mathbb{E}(\mathbf{g}_j^\top (\mathbf{I} - \mathbf{S})^\top \mathbf{\Gamma}^{-1} \mathbf{S} \boldsymbol{\varepsilon})^2 &\leq N_n^{-1} \|\mathbf{g}_j^\top (\mathbf{I} - \mathbf{S})^\top \mathbf{\Gamma}^{-1}\|^2 \mathbb{E}\|\mathbf{S} \boldsymbol{\varepsilon}\|^2 = \mathcal{O}(h^3 \log N_n), \\ N_n^{-1} \mathbb{E}(\boldsymbol{\phi}_j^\top \mathbf{S}^\top \mathbf{\Gamma}^{-1} \boldsymbol{\varepsilon})^2 &\leq N_n^{-1} \|\boldsymbol{\phi}_j^\top \mathbf{S}^\top\|^2 \|\mathbf{\Gamma}^{-1}\|_2^2 = \mathcal{O}(N_n^{-1} h^{-2} (\log N_n)^2), \\ N_n^{-1} \mathbb{E}(\mathbf{g}_j^\top (\mathbf{I} - \mathbf{S})^\top \mathbf{\Gamma}^{-1} \boldsymbol{\varepsilon})^2 &\leq N_n^{-1} \|\mathbf{g}_j^\top (\mathbf{I} - \mathbf{S})^\top\|^2 \|\mathbf{\Gamma}^{-1}\|_2^2 = \mathcal{O}(h^4).\end{aligned}$$

By Lemma 1 and Assumption (C.2), $\|\mathbf{S} \boldsymbol{\phi}_j\| = \mathcal{O}(N_n^{1/2} N_n^{-1} h^{-1} N_n^{1/2} \log N_n) = \mathcal{O}(h^{-1} \log N_n)$. Thus, for (I_{16}) ,

$$N_n^{-1} \mathbb{E}(\boldsymbol{\phi}_j^\top \mathbf{S}^\top \mathbf{\Gamma}^{-1} \mathbf{S} \boldsymbol{\varepsilon})^2 \leq N_n^{-1} \|\boldsymbol{\phi}_j^\top \mathbf{S}^\top\|^2 \|\mathbf{\Gamma}^{-1}\|_2^2 \mathbb{E}\|\mathbf{S} \boldsymbol{\varepsilon}\|^2 = \mathcal{O}(N_n^{-1} h^{-3} (\log N_n)^3).$$

Similarly, for (I_2) , we obtain

$$\begin{aligned}N_n^{-1/2} \mathbf{g}_j^\top (\mathbf{I} - \mathbf{S})^\top \mathbf{\Gamma}^{-1} (\mathbf{I} - \mathbf{S}) \mathbf{f} &= \mathcal{O}(N_n^{-1/2} N_n^{1/2} h^2 N_n^{1/2} h^2) = \mathcal{O}(N_n^{1/2} h^4), \\ N_n^{-1/2} \boldsymbol{\phi}_j^\top \mathbf{\Gamma}^{-1} (\mathbf{I} - \mathbf{S}) \mathbf{f} &= \mathcal{O}(N_n^{-1/2} N_n^{1/2} \log N_n h^2) = \mathcal{O}(h^2 \log N_n), \\ N_n^{-1/2} \boldsymbol{\phi}_j^\top \mathbf{S}^\top \mathbf{\Gamma}^{-1} (\mathbf{I} - \mathbf{S}) \mathbf{f} &= \mathcal{O}(N_n^{-1/2} h^{-1} \log N_n N_n^{1/2} h^2) = \mathcal{O}(h \log N_n).\end{aligned}$$

Thus, $N_n^{-1/2} \ell'(\boldsymbol{\beta}) \xrightarrow{D} N(\mathbf{0}, \mathbf{\Pi})$. The k th column of $-\ell''(\boldsymbol{\beta}, \boldsymbol{\theta})$ is $\mathbf{X}^\top (\mathbf{I} - \mathbf{S})^\top \mathbf{\Gamma}^k (\mathbf{I} - \mathbf{S}) (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) = \mathbf{X}^\top (\mathbf{I} - \mathbf{S})^\top \mathbf{\Gamma}^k (\mathbf{I} - \mathbf{S}) (\mathbf{f} + \boldsymbol{\varepsilon})$. The same argument can be used to show $N_n^{-1} \ell''(\boldsymbol{\beta}, \boldsymbol{\theta}) \xrightarrow{p} \mathbf{0}$.

Similarly, we can show

$$N_n^{-1} \ell''(\boldsymbol{\beta}) = -N_n^{-1} \mathbf{X}^\top (\mathbf{I} - \mathbf{S})^\top \mathbf{\Gamma}^{-1} (\mathbf{I} - \mathbf{S}) \mathbf{X} \xrightarrow{p} -N_n^{-1} \boldsymbol{\Phi}^\top \mathbf{\Gamma}^{-1} \boldsymbol{\Phi} = -\mathbf{\Pi}.$$

Therefore, $N_n^{-1} \ell''(\boldsymbol{\beta}, \boldsymbol{\beta}) \xrightarrow{p} -\mathbf{\Pi}$.

By Lemmas 2–3,

$$\|(\mathbf{I} - \mathbf{S}) \mathbf{f}\|^2 = N_n \mathcal{O}(h^4) = \mathcal{O}(N_n h^4), \quad \|\mathbf{S} \boldsymbol{\varepsilon}\|^2 = N_n \mathcal{O}_p\left(\frac{\log N_n}{N_n h}\right) = \mathcal{O}_p\left(\frac{\log N_n}{h}\right).$$

Since the k th element of $-2\ell'(\boldsymbol{\theta})$ is $\text{tr}(\mathbf{\Gamma}^{-1} \mathbf{\Gamma}^k) + (\mathbf{f} + \boldsymbol{\varepsilon})^\top (\mathbf{I} - \mathbf{S})^\top \mathbf{\Gamma}^k (\mathbf{I} - \mathbf{S}) (\mathbf{f} + \boldsymbol{\varepsilon})$, we can further show that

$$\begin{aligned}N_n^{-1/2} \mathbf{f}^\top (\mathbf{I} - \mathbf{S})^\top \mathbf{\Gamma}^k (\mathbf{I} - \mathbf{S}) \mathbf{f} &\leq N_n^{-1/2} \|\mathbf{\Gamma}^k\|_2 \|(\mathbf{I} - \mathbf{S}) \mathbf{f}\|^2 = \mathcal{O}(N_n^{1/2} h^4), \\ N_n^{-1/2} \mathbf{f}^\top (\mathbf{I} - \mathbf{S})^\top \mathbf{\Gamma}^k \boldsymbol{\varepsilon} &\leq N_n^{-1/2} \mathcal{O}(h^2) \mathbf{1}_{N_n}^\top \boldsymbol{\varepsilon} \xrightarrow{p} 0, \\ N_n^{-1/2} |\mathbf{f}^\top (\mathbf{I} - \mathbf{S})^\top \mathbf{\Gamma}^k \mathbf{S} \boldsymbol{\varepsilon}| &\leq N_n^{-1/2} \mathcal{O}_p(N_n^{1/2} h^3 \log N_n), \\ \mathbb{E} |N_n^{-1/2} \boldsymbol{\varepsilon}^\top \mathbf{S}^\top \mathbf{\Gamma}^k \mathbf{S} \boldsymbol{\varepsilon}| &\leq N_n^{-1/2} \|\mathbf{\Gamma}^k\|_2 \mathbb{E}\|\mathbf{S} \boldsymbol{\varepsilon}\|^2 = \mathcal{O}_p\left(\frac{\log N_n}{N_n^{1/2} h}\right).\end{aligned}$$

By Lemma 3, $N_n^{-1/2} \boldsymbol{\varepsilon}^\top \mathbf{S}^\top \mathbf{\Gamma}^k \boldsymbol{\varepsilon} = N_n^{-1/2} \mathbf{1}_{N_n}^\top \boldsymbol{\varepsilon} \mathcal{O}_p(1) \xrightarrow{p} 0$. Therefore, $N_n^{-1/2} \boldsymbol{\varepsilon}^\top (\mathbf{I} - \mathbf{S})^\top \mathbf{\Gamma}^k (\mathbf{I} - \mathbf{S}) \boldsymbol{\varepsilon} \xrightarrow{D} N_n^{-1/2} \boldsymbol{\varepsilon}^\top \mathbf{\Gamma}^k \boldsymbol{\varepsilon}$. Thus, $N_n^{-1/2} \ell'(\boldsymbol{\theta}) \xrightarrow{D} N(\mathbf{0}, \mathcal{I}_0(\boldsymbol{\theta}))$, and similar argument can be applied to show that $N_n^{-1} \ell''(\boldsymbol{\theta}, \boldsymbol{\theta}) \xrightarrow{p} -\mathcal{I}_0(\boldsymbol{\theta})$.

Next, we show the consistency and asymptotic normality of parameter estimates. To establish $\|\hat{\boldsymbol{\eta}} - \boldsymbol{\eta}_0\| = \mathcal{O}_p(N_n^{-1/2})$, it suffices to show that, for a given constant $\epsilon > 0$, there is a constant C such that, for a sufficiently large n ,

$$P\left\{\sup_{\|\mathbf{u}\|=C} \ell(\boldsymbol{\eta}_0 + N_n^{-1/2} \mathbf{u}) < \ell(\boldsymbol{\eta}_0)\right\} \geq 1 - \epsilon, \quad (15)$$

where $\mathbf{u} \in \mathbb{R}^{p+q}$. By Taylor's expansion, we obtain

$$\ell(\boldsymbol{\eta}_0 + N_n^{-1/2} \mathbf{u}) - \ell(\boldsymbol{\eta}_0) = N_n^{-1/2} \ell'(\boldsymbol{\eta}_0)^\top \mathbf{u} - (1/2) N_n^{-1} \mathbf{u}^\top \ell''(\boldsymbol{\eta}_0) \mathbf{u} \{1 + o_p(1)\}. \quad (16)$$

By (14a)–(14b), we have $N_n^{-1/2} \ell'(\eta_0) = \mathcal{O}_p(1)$ and $N_n^{-1} \ell''(\eta_0) = \mathcal{O}_p(1)$. Therefore, for a sufficiently large C , the second term dominates the first term in (16), and therefore, (15) holds.

To further establish the asymptotic normality of $\hat{\eta}$, it can be shown that $\hat{\eta} = (\hat{\beta}^\top, \hat{\theta}^\top)^\top$ satisfies

$$\mathbf{0} = \ell'(\eta)|_{\eta=\hat{\eta}} = \ell'(\eta_0) + \{\ell''(\eta_0) + o_p(1)\}(\hat{\eta} - \eta_0).$$

Together with (14a)–(14b), Theorem 1 holds. \square

Proof of Theorem 2. Write $\hat{F}(t) - F(t) = \omega(t) \{\mathbf{f} + \boldsymbol{\varepsilon} - \mathbf{X}(\hat{\beta} - \beta)\} - F(t) = (II_1) + (II_2) - (II_3)$, where $(II_1) = \omega(t)\mathbf{f} - F(t)$, $(II_2) = \omega(t)\boldsymbol{\varepsilon}$ and $(II_3) = \omega(t)\mathbf{X}(\hat{\beta} - \beta)$.

For (II_1) , a Taylor's expansion yields

$$\begin{aligned} f(t_i) &= f(t) + f'(t)(t_i - t) + 1/2 f''(t)(t_i - t)^2 + 1/6 f^{(3)}(\xi_i)(t_i - t)^3 \\ &= \left(1, \frac{t_i - t}{h}\right) \mathbf{F}(t) + 1/2 f''(t)(t_i - t)^2 + 1/6 f^{(3)}(\xi_i)(t_i - t)^3, \end{aligned}$$

where ξ_i is between t and t_i . Therefore, $(II_1) = (1/2)\omega(t)\mathbf{d}_2 f''(t) + \frac{1}{6}\omega(t)\mathbf{d}_{3,\xi}$, where $\mathbf{d}_2 = ((t_1 - t)^2, \dots, (t_{N_n} - t)^2)^\top$ and $\mathbf{d}_{3,\xi} = (f^{(3)}(\xi_1)(t_1 - t)^3, \dots, f^{(3)}(\xi_{N_n})(t_{N_n} - t)^3)^\top$. Given $t \in (0, 1)$, by Lemma 1,

$$\omega(t)\mathbf{d}_2 f''(t) = \begin{pmatrix} v_{2,t}^2 - v_{1,t}v_{3,t} \\ v_{0,t}v_{3,t} - v_{1,t}v_{2,t} \end{pmatrix} \frac{h^2 f''(t)}{v_{0,t}v_{2,t} - v_{1,t}^2} = h^2 \begin{pmatrix} \mu_2 f''(t) \\ 0 \end{pmatrix} + o(h^2).$$

Moreover, we have

$$|\omega(t)\mathbf{d}_{3,\xi}| \leq \max_{x \in \mathbb{R}} \frac{4|f^{(3)}(x)|h^3}{v_{0,t}v_{2,t} - v_{1,t}^2} \left\{ \left| \begin{pmatrix} v_{2,t}v_{3,t} - v_{1,t}v_{4,t} \\ -v_{1,t}v_{3,t} + v_{0,t}v_{4,t} \end{pmatrix} \right| \right\} = \mathcal{O}(h^3).$$

Therefore, $(II_1) = h^2 \begin{pmatrix} \mu_2 f''(t) \\ 0 \end{pmatrix} + o(h^2)$.

For (II_2) , let $A(\boldsymbol{\varepsilon}) = \sum_{i=1}^{N_n} K_h(t_i - t)\varepsilon_i$ and $B(\boldsymbol{\varepsilon}) = \sum_{i=1}^{N_n} K_h(t_i - t)\frac{t_i - t}{h}\varepsilon_i$, we have

$$\omega(t)\boldsymbol{\varepsilon} = \frac{1}{v_{0,t}v_{2,t} - v_{1,t}^2} \begin{pmatrix} v_{2,t}A(\boldsymbol{\varepsilon}) - v_{1,t}B(\boldsymbol{\varepsilon}) \\ -v_{1,t}A(\boldsymbol{\varepsilon}) + v_{0,t}B(\boldsymbol{\varepsilon}) \end{pmatrix}.$$

For $t \in (0, 1)$, by Lemma 1, we have $\frac{N_n v_{0,t}}{v_{0,t}v_{2,t} - v_{1,t}^2} \rightarrow \mu_2^{-1}q(t)^{-1}$, $\frac{N_n v_{1,t}}{v_{0,t}v_{2,t} - v_{1,t}^2} \rightarrow 0$ and $\frac{N_n v_{2,t}}{v_{0,t}v_{2,t} - v_{1,t}^2} \rightarrow q(t)^{-1}$. Since $\boldsymbol{\varepsilon}$ is a Gaussian process, by (C.9) and Slutsky's Theorem,

$$(N_n h)^{1/2} \omega(t)\boldsymbol{\varepsilon} \xrightarrow{D} N\left(0, \begin{pmatrix} 1 & 0 \\ 0 & \mu_2^{-1} \end{pmatrix} \Delta_t \begin{pmatrix} 1 & 0 \\ 0 & \mu_2^{-1} \end{pmatrix}\right).$$

For (II_3) , by Lemma 6, $\omega(t)\mathbf{X} = \mathcal{O}(1)$. By Theorem 1, we have $\hat{\beta} - \beta = \mathcal{O}_p(N_n^{-1/2})$, and therefore, $(II_3) = \mathcal{O}_p(N_n^{-1/2})$. \square

Proof of Theorem 3. Let $\omega_1^{(-i)}(t) = (1, 0) \left[\left\{ D_t^{(-i)} \right\}^\top K_t^{(-i)} D_t^{(-i)} \right]^{-1} \left\{ D_t^{(-i)} \right\}^\top K_t^{(-i)}$, where $D_t^{(-i)}$ is the matrix of D_t with i th row deleted, and $K_t^{(-i)}$ is the matrix K_t with both i th row and column deleted. In addition, we let $\mathbf{y}^{*(-i)}$ denote the vector of response variables with the i th entry left out. Straightforward calculation reveals

$$\hat{f}^{(-i)}(t_i) = \omega_1^{(-i)}(t_i) \mathbf{y}^{*(-i)} = \frac{\begin{pmatrix} v_{2,t_i}^{(-i)} - v_{1,t_i}^{(-i)} \end{pmatrix}}{v_{2,t_i}^{(-i)} v_{0,t_i}^{(-i)} - (v_{1,t_i}^{(-i)})^2} \left\{ D_t^{(-i)} \right\}^\top K_{t_i}^{(-i)} \mathbf{y}^{*(-i)} = \sum_{j \neq i} a_j^{(-i)}(t_i) y_j^*,$$

where $a_j^{(-i)}(t_i) = \frac{\begin{pmatrix} v_{2,t_i}^{(-i)} - v_{1,t_i}^{(-i)} \end{pmatrix}}{v_{2,t_i}^{(-i)} v_{0,t_i}^{(-i)} - (v_{1,t_i}^{(-i)})^2} \left\{ D_t^{(-i)} \right\}^\top K_{t_i}^{(-i)}$, $v_{0,t_i}^{(-i)} = \sum_{j \neq i} K_h(t_j - t_i) = v_{0,t_i} - K_h(0)$, $v_{1,t_i}^{(-i)} = \sum_{j \neq i} K_h(t_j - t_i) \frac{t_j - t_i}{h} = v_{1,t_i}$ and $v_{2,t_i}^{(-i)} = \sum_{j \neq i} K_h(t_j - t_i) \left(\frac{t_j - t_i}{h} \right)^2 = v_{2,t_i}$.

By Lemma 1,

$$\begin{aligned} a_j^{(-i)}(t_i) &= \frac{v_{2,t_i}^{(-i)} K_h(t_j - t_i)}{v_{2,t_i}^{(-i)} v_{0,t_i}^{(-i)} - (v_{1,t_i}^{(-i)})^2} - \frac{v_{1,t_i}^{(-i)} K_h(t_j - t_i)(t_j - t_i)/h}{v_{2,t_i}^{(-i)} v_{0,t_i}^{(-i)} - (v_{1,t_i}^{(-i)})^2} = \frac{v_{2,t_i} K_h(t_j - t_i)}{\{v_{0,t_i} - K_h(0)\} v_{2,t_i} - v_{1,t_i}^2} - \frac{v_{1,t_i} K_h(t_j - t_i)(t_j - t_i)/h}{\{v_{0,t_i} - K_h(0)\} v_{2,t_i} - v_{1,t_i}^2} \\ &= \frac{K\left(\frac{t_j - t_i}{h}\right) - \frac{\mu_{1,t_i}}{\mu_{2,t_i}} K\left(\frac{t_j - t_i}{h}\right) \left(\frac{t_j - t_i}{h}\right)}{N_n h q(t_i) (\mu_{0,t_i} \mu_{2,t_i} - \mu_{1,t_i}^2) / \mu_{2,t_i} - K(0)} + \mathcal{O}(N_n^{-1} + N_n^{-1} \zeta_n h^{-2}). \end{aligned}$$

Thus, for the leave-one-out cross-validation (CV) score function,

$$E\{CV(h)\} = E\left[\frac{1}{N_n} \sum_{i=1}^{N_n} \{f(t_i) + \varepsilon_i - \hat{f}^{(-i)}(t_i)\}^2\right] = \frac{1}{N_n} \sum_{i=1}^{N_n} E\{f(t_i) - \hat{f}^{(-i)}(t_i)\}^2 + \frac{1}{N_n} \sum_{i=1}^{N_n} \text{Var}(Y_i) - \frac{2}{N_n} \sum_{i=1}^{N_n} \text{Cov}\{\hat{f}^{(-i)}(t_i), \varepsilon_i\}.$$

Denote $A(h) = \frac{2}{N_n} \sum_{i=1}^{N_n} \text{Cov}\{\hat{f}^{(-i)}(t_i), \varepsilon_i\}$, we have

$$A(h) = \frac{2}{N_n} \sum_{i=1}^{N_n} \sum_{j \neq i} \frac{K\left(\frac{t_j - t_i}{h}\right) \left(1 - \frac{\mu_{1,t_i}}{\mu_{2,t_i}} \left(\frac{t_j - t_i}{h}\right)\right)}{b(t_i) - K(0)} \text{Cov}(\varepsilon_i, \varepsilon_j) + o(N_n^{-1}),$$

since $\|\Gamma\|_\infty = O(1)$ as shown in the proof of Theorem 1 in Chu et al. [5]. Under the asymptotic framework (A.1), and since $K(\cdot)$ has a bounded first-order derivative at the origin, we obtain

$$\frac{2}{N_n} \sum_{i=1}^{N_n} \sum_{\substack{j \neq i \\ |t_j - t_i| \leq C_n}} \frac{K\left(\frac{t_j - t_i}{h}\right) \left(1 - \frac{\mu_{1,t_i}}{\mu_{2,t_i}} \left(\frac{t_j - t_i}{h}\right)\right)}{b(t_i) - K(0)} \text{Cov}(\varepsilon_i, \varepsilon_j) = \frac{K(0)}{b(t_i) - K(0)} \left(\frac{2}{N_n} \sum_{i=1}^{N_n} \sum_{\substack{j \neq i \\ B_n |t_j - t_i| < C_n}} \text{Cov}(\varepsilon_i, \varepsilon_j) \right) + o\left(\frac{1}{N_n h}\right).$$

Note that

$$\sum_{\substack{j \neq i \\ |t_j - t_i| > C_n}} \text{Cov}(\varepsilon_i, \varepsilon_j) = \sum_{m' = \lfloor \frac{B_n C_n}{b} \rfloor} \mathcal{O}\left(\frac{b}{B_n \zeta_n}\right) \max_{mb \leq |u_2| \leq (m+1)b} \gamma_1(|u|) \leq \mathcal{O}\left(\frac{b}{B_n \zeta_n}\right) \int_{B_n C_n}^{\infty} \gamma_1(u) du \rightarrow 0,$$

so we have

$$\frac{2}{N_n} \sum_{i=1}^{N_n} \sum_{|t_j - t_i| > C_n} \frac{K\left(\frac{t_j - t_i}{h}\right) \left(1 - \frac{\mu_{1,t_i}}{\mu_{2,t_i}} \left(\frac{t_j - t_i}{h}\right)\right)}{b(t_i) - K(0)} \text{Cov}(\varepsilon_i, \varepsilon_j) = o\left(\frac{1}{N_n h}\right).$$

Therefore, $A(h) = K(0) \left(\frac{2}{N_n} \sum_{i=1}^{N_n} \sum_{\substack{j \neq i \\ |t_j - t_i| < C_n}} \frac{\text{Cov}(\varepsilon_i, \varepsilon_j)}{b(t_i) - K(0)} \right) + o\left(\frac{1}{N_n h}\right)$. Thus, the desired results are shown. \square

Appendix B. Supplementary material

Supplementary material related to this article can be found online at <https://doi.org/10.1016/j.jmva.2021.104735>.

References

- [1] D. Al-Sulami, Z. Jiang, Z. Lu, J. Zhu, On a semiparametric data-driven nonlinear model with penalized spatio-temporal lag interactions, *J. Time Series Anal.* 40 (2019) 327–342.
- [2] N.S. Altman, Kernel smoothing of data with correlated errors, *J. Am. Stat. Assoc.* 85 (1990) 749–759.
- [3] S. Bandyopadhyay, C. Jentsch, S. Subba Rao, A spectral domain test for stationarity of spatio-temporal data, *J. Time Series Anal.* 38 (2017) 326–351.
- [4] S. Bandyopadhyay, S.S. Rao, A test for stationarity for irregularly spaced spatial data, *J. R. Stat. Soc. Ser. B Stat. Methodol.* 79 (2017) 95–123.
- [5] T. Chu, J. Liu, J. Zhu, H. Wang, Spatio-temporal expanding distance asymptotic framework for locally stationary processes, *Sankhya A* (2020) 1–25.
- [6] T. Chu, J. Zhu, H. Wang, Penalized maximum likelihood estimation and variable selection in geostatistics, *Ann. Statist.* 39 (2011) 2607–2625.
- [7] T. Chu, J. Zhu, H. Wang, Semiparametric modeling with nonseparable and nonstationary spatio-temporal covariance functions and its inference, *Statist. Sinica* 29 (2019) 1233–1252.
- [8] N. Cressie, *Statistics for Spatial Data*, revised edition, Wiley, New York, 1993.
- [9] N. Cressie, H.-C. Huang, Classes of nonseparable spatio-temporal stationary covariance functions, *J. Amer. Statist. Assoc.* 94 (1999) 1330–1340.
- [10] N. Cressie, S. Lahiri, The asymptotic distribution of REML estimators, *J. Multivariate Anal.* 45 (1993) 217–233.
- [11] A. Datta, S. Banerjee, A.O. Finley, N.A. Hamm, M. Schaap, Nonseparable dynamic nearest neighbor gaussian process models for large spatio-temporal data with an application to particulate matter analysis, *Ann. Appl. Stat.* 10 (2016) 1286.
- [12] K. De Brabanter, J. De Brabanter, J.A.K. Suykens, B. De Moor, Kernel regression in the presence of correlated errors, *J. Mach. Learn. Res.* 12 (2011) 1955–1976.
- [13] J. Fan, I. Gijbels, *Local Polynomial Modelling and its Applications*, Chapman and Hall/CRC, London, 1996.
- [14] J. Fan, T. Huang, Profile likelihood inferences on semiparametric varying-coefficient partially linear models, *Bernoulli* 11 (2005) 1031–1057.
- [15] M. Francisco-Fernandez, J.D. Opsomer, Smoothing parameter selection methods for nonparametric regression with spatially correlated errors, *Canad. J. Statist.* 33 (2005) 279–295.
- [16] M. Fuentes, L. Chen, J.M. Davis, A class of nonseparable and nonstationary spatial temporal covariance functions, *Environmetrics* 19 (2008) 487–507.
- [17] J. Gao, H. Liang, Statistical inference in single-index and partially nonlinear models, *Ann. Inst. Statist. Math.* 49 (1997) 493–517.
- [18] J. Gao, Z. Lu, D. Tjøstheim, Estimation in semiparametric spatial regression, *Ann. Statist.* 34 (2006) 1395–1435.
- [19] T. Gneiting, Nonseparable, stationary covariance functions for space-time data, *J. Amer. Statist. Assoc.* 97 (2002) 590–600.
- [20] G.H. Golub, M. Heath, G. Wahba, Generalized cross-validation as a method for choosing a good ridge parameter, *Technometrics* 21 (1979) 215–223.

- [21] P. Hall, P. Patil, Properties of nonparametric estimators of autocovariance for stationary random fields, *Probab. Theory Related Fields* 99 (1994) 399–424.
- [22] W. Hä, H. Liang, J.T. Gao, *Partially Linear Models*, Springer, New York, 2000.
- [23] S. Lahiri, Central limit theorems for weighted sums of a spatial process under a class of stochastic and fixed designs, *Sankhya A* 65 (2003a) 356–388.
- [24] S. Lahiri, *Resampling Methods for Dependent Data*, Springer, New York, 2003b.
- [25] K. Lake, J. Zhu, H. Wang, J. Volckens, K.A. Koehler, Effects of data sparsity and spatiotemporal variability on hazard maps of workplace noise, *J. Occup. Environ. Hygiene* 12 (2015) 256–265.
- [26] H. Liang, W. Härdle, R.J. Carroll, Estimation in a semiparametric partially linear errors-in-variables model, *Ann. Statist.* 27 (1999) 1519–1535.
- [27] H. Liang, R. Li, Variable selection for partially linear models with measurement errors, *J. Amer. Statist. Assoc.* 104 (2009) 234–248.
- [28] W.-L. Loh, Fixed-domain asymptotics for a subclass of Matérn-type gaussian random fields, *Ann. Statist.* 33 (2005) 2344–2394.
- [29] Z. Lu, D.J. Steinskog, D. Tjøstheim, Q. Yao, Adaptively varying coefficient spatiotemporal models, *J. R. Stat. Soc. Ser. B Stat. Methodol.* 71 (2009) 859–880.
- [30] Z. Lu, D. Tjøstheim, Nonparametric estimation of probability density functions for irregularly observed spatial data, *J. Amer. Statist. Assoc.* 109 (2014) 1546–1564.
- [31] G. Ludwig, T. Chu, J. Zhu, H. Wang, K. Koehler, Static and roving sensor data fusion for spatio-temporal hazard mapping with application to occupational exposure assessment, *Ann. Appl. Stat.* 11 (2017) 139–160.
- [32] K.V. Mardia, R.J. Marshall, Maximum likelihood estimation of models for residual covariance in spatial regression, *Biometrika* 71 (1984) 135–146.
- [33] J. Opsomer, Y. Wang, Y. Yang, Nonparametric regression with correlated errors, *Statist. Sci.* 16 (2001) 134–153.
- [34] E. Porcu, A. Alegria, R. Furrer, Modeling temporally evolving and spatially globally dependent data, *Internat. Statist. Rev.* 86 (2018) 344–377.
- [35] A. Rodrigues, P.J. Diggle, A class of convolution-based models for spatio-temporal processes with non-separable covariance structure, *Scand. J. Stat.* 37 (2010) 553–567.
- [36] D. Ruppert, M. Wand, R. Carroll, *Semiparametric Regression*, Cambridge University Press, New York, 2003.
- [37] P. Speckman, Kernel smoothing in partial linear models, *J. R. Stat. Soc. Ser. B Stat. Methodol.* 50 (1988) 413–436.
- [38] M.L. Stein, *Interpolation of Spatial Data: Some Theory for Kriging*, Springer, New York, 1999.
- [39] M.L. Stein, Space-time covariance functions, *J. Amer. Statist. Assoc.* 100 (2005) 310–321.
- [40] C.J. Stone, Optimal global rates of convergence for nonparametric regression, *Ann. Statist.* 10 (1982) 1040–1053.
- [41] J.R. Stroud, P. Müller, B. Sansó, Dynamic models for spatiotemporal data, *J. R. Stat. Soc. Ser. B Stat. Methodol.* 63 (2001) 673–689.
- [42] Y. Sun, H. Yan, W. Zhang, Z. Lu, A semiparametric spatial dynamic model, *Ann. Statist.* 42 (2014) 700–727.
- [43] M. Vogt, O. Linton, Nonparametric estimation of a periodic sequence in the presence of a smooth trend, *Biometrika* 101 (2014) 121–140.
- [44] L. Wasserman, *All of Nonparametric Statistics*, Springer, New York, 2010.
- [45] Z. Ying, Maximum likelihood estimation of parameters under a spatial sampling scheme, *Ann. Statist.* 21 (1993) 1567–1590.
- [46] H. Zhang, Inconsistent estimation and asymptotically equivalent interpolations in model-based geostatistics, *J. Amer. Statist. Assoc.* 99 (2004) 250–261.