

# An Efficient Implementation for Spatial-Temporal Gaussian Process Regression and Its Applications <sup>★</sup>

Junpeng Zhang <sup>a</sup>, Yue Ju <sup>a</sup>, Biqiang Mu <sup>b</sup>, Renxin Zhong <sup>c</sup>, Tianshi Chen <sup>a</sup>

<sup>a</sup>*School of Data Science and Shenzhen Research Institute of Big Data, The Chinese University of Hong Kong, Shenzhen 518172, China*

<sup>b</sup>*Key Laboratory of Systems and Control, Institute of Systems Science, Academy of Mathematics and Systems Science, Chinese Academy of Sciences, Beijing 100190, China*

<sup>c</sup>*School of Intelligent Systems Engineering, Sun Yat-sen University, Guangzhou, P. R. China*

---

## Abstract

Spatial-temporal Gaussian process regression is a popular method for spatial-temporal data modeling. Its state-of-art implementation is based on the state-space model realization of the spatial-temporal Gaussian process and its corresponding Kalman filter and smoother, and has computational complexity  $\mathcal{O}(NM^3)$ , where  $N$  and  $M$  are the number of time instants and spatial input locations, respectively, and thus can only be applied to data with large  $N$  but relatively small  $M$ . In this paper, our primary goal is to show that by exploring the Kronecker structure of the state-space model realization of the spatial-temporal Gaussian process, it is possible to further reduce the computational complexity to  $\mathcal{O}(M^3 + NM^2)$  and thus the proposed implementation can be applied to data with large  $N$  and moderately large  $M$ . The proposed implementation is illustrated over applications in weather data prediction and spatially-distributed system identification. Our secondary goal is to design a kernel for both the Colorado precipitation data and the GHCN temperature data, such that while having more efficient implementation, better prediction performance can also be achieved than the state-of-art result.

*Key words:* Large scale spatial-temporal data; Gaussian process regression; Kalman filter and smoother.

---

## 1 Introduction

Gaussian process regression is a popular method in statistical data modeling and analysis, closely related with the kernel method, e.g., [16], and the kernel-based regularization method in system identification, e.g., [24], and has wide applications in many fields such as machine learning, signal processing, and automatic control,

e.g., [23, 24, 29]. In contrast with the common parametric modeling methods in system identification, e.g., the prediction error/maximum likelihood method [19], its advantage lies in that, first, its model structure is determined by its covariance function (also called the kernel function), which incorporates the prior knowledge of the underlying function/system to be identified into the estimation procedure; second, its model complexity is governed by the parameter (called the hyper-parameter) used to parameterize the covariance function, and can be tuned in a continuous way. To apply Gaussian process regression methods, there are several issues that should be addressed, including the kernel design, e.g., [9, 10, 14, 30], the hyper-parameter estimation, e.g., [24, 29], and the efficient implementation, e.g., [4, 10, 25]. Gaussian process regression has been used widely in dealing with the spatial-temporal data in many areas, such as climate science, social science, kriging, signal processing and physical inverse problems, e.g., [3, 13, 22, 26]. When dealing with the spatial-temporal data, the Gaussian process has two inputs: the locations and the time instants, and thus is often referred to as the spatial-temporal Gaussian process, e.g., [7, 26, 27].

---

<sup>★</sup> A preliminary version of this work [31] was presented in the 39th Chinese Control Conference (CCC), 2020. Corresponding author Tianshi Chen. This work was supported by the Shenzhen Science and Technology Innovation Council under contract No. Ji-20170189 (JCYJ20170411102101881), the Robotic Discipline Development Fund (2016-1418) from Shenzhen Government, the general project funded by NSFC under contract No. 61773329, the Thousand Youth Talents Plan funded by the central government of China.

*Email addresses:* junpengzhang@link.cuhk.edu.cn (Junpeng Zhang), yueju@link.cuhk.edu.cn (Yue Ju), bqmu@amss.ac.cn (Biqiang Mu), zhrenxin@mail.sysu.edu.cn (Renxin Zhong), tschen@cuhk.edu.cn (Tianshi Chen).

For large scale spatial-temporal data, the aforementioned issues often become more involved. The current practice is to design a separable spatial-temporal kernel, which is a product of a spatial kernel and a temporal kernel, e.g., [7, 26, 27], whose design should be based on the prior knowledge on the underlying function to be identified. For a designed separable spatial-temporal kernel, many methods can be used for the hyper-parameter estimation, such as the empirical Bayes (EB) method (also called the marginal likelihood maximization (MLM) method), the Stein’s unbiased risk estimate (SURE) minimization method, the generalized cross validation (GCV) method, e.g., [24, 29]. The straightforward implementation of the hyper-parameter estimation and the following estimation and prediction step has computational complexity  $\mathcal{O}(N^3M^3)$ , where  $N$  and  $M$  are the numbers of temporal and spatial data, respectively, and thus is too expensive to be applied to large scale data. To reduce the computational complexity, it has been tried to first explore the structure of the temporal kernel, then derive a state-space model realization of the temporal Gaussian process in different ways, and finally convert the hyper-parameter estimation, function estimation and prediction to Kalman filtering, smoothing and prediction problems, e.g., [7, 17, 26, 27]. Such implementation has computational complexity  $\mathcal{O}(NM^3)$  and thus is inefficient to be applied to spatial-temporal data with moderately large or large  $M$ .

In this paper, we focus on the following two issues. Our primary focus is on the issue of how to further reduce the computational complexity such that the Gaussian process regression can be applied to spatial-temporal data with large  $N$  and moderately large  $M$ . To tackle this problem, we first find that the state-space model realization of the spatial-temporal Gaussian process has a Kronecker structure and then by exploring this structure, we propose transformations for the original state-space model and then derive a new state-space model realization of the spatial-temporal Gaussian process. Finally, the Kalman filter, smoother and predictor are applied to handle the hyper-parameter estimation, function estimation and prediction, respectively. The proposed implementation is illustrated over applications in weather data prediction including the Colorado precipitation data considered in [27] and the Global Historical Climatology Network (GHCN) temperature data in [21, 31], and spatially-distributed system identification, e.g., [18]. Our secondary focus is to design a kernel for the two weather data sets, such that while having more efficient implementation, better prediction performance can also be achieved than the kernel proposed in [27]. To this purpose, the designed kernel should have state-space model realizations and also incorporate the prior knowledge that both data sets are not strictly periodic but with slight temporal variation.

In contrast with the state-of-art result [27], this paper has the following contributions:

- 1) a more efficient implementation algorithm with computational complexity  $\mathcal{O}(M^3 + NM^2)$  is proposed for the hyper-parameter estimation, the spatial-temporal Gaussian process regression and prediction, while the one in [27] has a computational complexity  $\mathcal{O}(NM^3)$  and did not consider the efficient implementation of hyper-parameter estimation;
- 2) a kernel is designed for the Colorado precipitation data and the GHCN temperature data and shown to give better prediction performance than the one in [27].

Finally, in contrast with the preliminary version [31] of this paper, we have included new theoretical results including Propositions 1 to 3, and Theorem 1, designed a new kernel that gives better prediction performance for both the Colorado precipitation data and the GHCN temperature data, illustrated the implementation over a new application in spatially-distributed system identification, and included more implementation details, e.g., the derivation of the discrete-time state-space model realization of the temporal kernel, the treatment of the missing data and the selection of the starting points.

The remaining parts of this paper are organized as follows. In Section 2, we first introduce some preliminary materials and then the problem statement. In Section 3, we propose an implementation with computational complexity  $\mathcal{O}(M^3 + NM^2)$ . In Section 4, we test the proposed implementation over applications in weather data prediction and spatially-distributed system identification, where in Section 4.2, we design a kernel and show its better prediction capability over the one in [27] for both the Colorado precipitation data and the GHCN temperature data. In Section 5, we give the conclusion of this paper. All proofs of theorems and propositions are included in Appendix A.

## 2 Preliminary and Problem Statement

In this section, we first introduce some preliminary materials and then the problem statement of this paper.

### 2.1 Spatial-temporal Function Estimation

In this paper, we consider the spatial-temporal function estimation problem described by

$$y_{i,j} = f(p_i, t_j) + v_{i,j}, \quad (1)$$

$$i = 1, \dots, M, \quad t_j = jT_s, \quad j = 1, \dots, N + N_T,$$

where  $p_i \in \mathbb{R}^\nu$  with  $\nu \in \mathbb{N}$  is the  $i$ th location,  $t_j \in \mathbb{R}_+ = \{x | x \geq 0, x \in \mathbb{R}\}$  is the  $j$ th time instant,  $f(p_i, t_j)$ ,  $v_{i,j} \in \mathbb{R}$  and  $y_{i,j} \in \mathbb{R}$  are the unknown spatial-temporal function value, the measurement noise and the measurement output at the  $i$ th location and the  $j$ th time instant, respectively,  $T_s > 0$  the sampling interval,  $M$  is

the number of locations, and  $N$  and  $N_T$  are the numbers of time instants for the function estimation and validation, respectively. The measurement noises  $v_{i,j}$  with  $i = 1, \dots, M, j = 1, \dots, N + N_T$  are assumed to be independently Gaussian distributed as follows

$$v_{i,j} \sim \mathcal{N}(0, \sigma^2). \quad (2)$$

We aim to estimate the function  $f : \mathbb{R}^m \times \mathbb{R}_+ \rightarrow \mathbb{R}$  based on the training data  $\{p_i, t_j, y_{i,j}\}_{i=1, j=1}^{M, N}$  such that it has as good prediction performance over the test data  $\{p_i, t_j, y_{i,j}\}_{i=1, j=N+1}^{M, N+N_T}$  as possible.

## 2.2 Gaussian Process Regression

Gaussian process regression models the spatial-temporal function  $f(p_i, t_j)$  as a spatial-temporal Gaussian process

$$f(p_i, t_j) \sim \mathcal{GP}(0, k(p_i, t_j, p_{i'}, t_{j'}; \alpha)), \quad (3)$$

$$k(p_i, t_j, p_{i'}, t_{j'}; \alpha) = k_s(p_i, p_{i'}; \alpha_s) k_t(t_j, t_{j'}; \alpha_t), \quad (4)$$

where  $i, i' = 1, \dots, M, j, j' = 1, \dots, N + N_T$ ,  $\mathcal{GP}$  represents a Gaussian process,  $k(p_i, t_j, p_{i'}, t_{j'}; \alpha)$  the covariance function (also called the kernel) with a separable structure in space and time, e.g. [7, 26, 27],  $k_s(p_i, p_{i'}; \alpha_s) : \mathbb{R}^\nu \times \mathbb{R}^\nu \rightarrow \mathbb{R}$  the spatial kernel,  $k_t(t_j, t_{j'}; \alpha_t) : \mathbb{R}_+ \times \mathbb{R}_+ \rightarrow \mathbb{R}$  the temporal kernel,  $\alpha = [\alpha_t^T, \alpha_s^T]^T \in \Omega \subset \mathbb{R}^d$  with  $d \in \mathbb{N}$ ,  $\alpha_s \in \mathbb{R}^{d_s}$  with  $d_s \in \mathbb{N}$  and  $\alpha_t \in \mathbb{R}^{d_t}$  with  $d_t \in \mathbb{N}$  the hyper-parameters of  $k, k_s$  and  $k_t$ , respectively, and  $d = d_t + d_s$ . It is assumed that for any  $i, i' = 1, \dots, M, j, j' = 1, \dots, N + N_T$ ,  $f(p_{i'}, t_{j'})$  is independent of  $v_{i,j}$ .

The kernel  $k(p_i, t_j, p_{i'}, t_{j'}; \alpha)$  determines the underlying model structure and its design for the two test data sets will be studied in Section 4.2. The hyper-parameter  $\alpha$  determines the model complexity and its estimation can be handled by many methods. Here, we consider the marginal likelihood maximization (MLM) method, the generalized cross validation (GCV) method, and the Stein's unbiased risk estimation (SURE) method, e.g., [24], which are listed below, respectively,

$$\hat{\alpha}^{MLM} = \arg \min_{\alpha \in \Omega} \left\{ \frac{NM}{2} \log(2\pi) + \frac{1}{2} \log |\Sigma(\alpha)| + \frac{1}{2} Y^T \Sigma^{-1}(\alpha) Y \right\}, \quad (5)$$

$$\hat{\alpha}^{GCV} = \arg \min_{\alpha \in \Omega} \left\{ \frac{S}{NM(1 - \delta/NM)^2} \right\}, \quad (6)$$

$$\hat{\alpha}^{SURE} = \arg \min_{\alpha \in \Omega} \{S + 2\sigma^2 \delta\}, \quad (7)$$

where  $\hat{\alpha}^{MLM}$ ,  $\hat{\alpha}^{GCV}$  and  $\hat{\alpha}^{SURE}$  denote the correspond-

ing hyper-parameter estimate,

$$\Sigma(\alpha) = K_t(\alpha_t) \otimes K_s(\alpha_s) + \sigma^2 I_{NM} \in \mathbb{R}^{NM \times NM}, \quad (8a)$$

$$[K_t(\alpha_t)]_{jj'} = k_t(t_j, t_{j'}; \alpha_t), \quad j, j' = 1, \dots, N, \quad (8b)$$

$$[K_s(\alpha_s)]_{ii'} = k_s(p_i, p_{i'}; \alpha_s), \quad i, i' = 1, \dots, M, \quad (8c)$$

$$y_j = [y_{1,j}, \dots, y_{M,j}]^T \in \mathbb{R}^M, \quad (8d)$$

$$Y = [y_1^T, \dots, y_N^T]^T \in \mathbb{R}^{NM}, \quad (8e)$$

$$\delta = \text{trace} \{ [K_t(\alpha_t) \otimes K_s(\alpha_s)] \Sigma(\alpha)^{-1} \}, \quad (8f)$$

$$S = \|\hat{Y} - Y\|_2^2, \quad (8g)$$

$$\hat{Y} = [K_t(\alpha_t) \otimes K_s(\alpha_s)] \Sigma(\alpha)^{-1} Y, \quad (8h)$$

$\otimes$  denotes the Kronecker product between two matrices,  $I_{NM} \in \mathbb{R}^{NM \times NM}$  an  $NM$ -dimensional identity matrix,  $[\cdot]_{jj'}$  the  $(j, j')$ th entry of a matrix,  $|\cdot|$  and  $\text{trace}(\cdot)$  the determinant and trace of a square matrix, respectively.

## 2.3 Problem Statement

To state the problem, it is worth to note the following two observations. Firstly, the state-of-art implementation in [27] designed a kernel such that the spatial-temporal Gaussian process has a state-space model realization and then convert the function estimation problem to a Kalman filtering and smoothing problem, and the implementation has computational complexity  $\mathcal{O}(NM^3)$ , e.g., [7, 26, 27] and thus can be applied to data with large  $N$  but relatively small  $M$ , e.g., the Colorado precipitation data with  $N = 1212$  and  $M = 367$  was studied in [27]. However, the implementation in [27] is still very expensive to apply for data with moderately large  $M$ , e.g., the GHCN temperature data with  $N = 6575$  and  $M = 3955$ . Secondly, the kernel designed in [27] does not give very good prediction performance for the Colorado precipitation data [17], indicating there is a room to design better kernels.

The above observations motivate us to tackle the following two problems in this paper:

- 1) to develop implementation with lower computational complexity in terms of  $M$  than the one in [27], which can be applied to data with large  $N$  and moderately large  $M$ , e.g., the GHCN temperature data;
- 2) to design a kernel for both the Colorado precipitation data and the GHCN temperature data that gives better prediction performance than the one in [27].

## 3 An Efficient Implementation

In this section, we propose a new implementation algorithm with computational complexity  $\mathcal{O}(M^3 + NM^2)$ , which can thus be applied to data with large  $N$  and moderately large  $M$ .

### 3.1 State-space Model Realization of Spatial-Temporal Gaussian Process

For convenience, we assume in this section that the temporal kernel  $k_t(t_j, t_{j'}; \alpha_t)$  is a stationary kernel and then with a slight abuse of the notation, we can denote it by

$$k_t(\tau; \alpha_t), \quad \tau = (j - j')T_s, \quad j, j' = 1, \dots, N + N_T. \quad (9)$$

Recall that the power spectral density (PSD) of a discrete-time kernel  $k_t(\tau; \alpha_t)$ , denoted as  $\Phi(\omega)$ , can be obtained by its discrete Fourier transform

$$\Phi(\omega) = \sum_{\tau=-\infty}^{+\infty} k_t(\tau; \alpha_t) e^{-i\omega\tau} \in \mathbb{R}_+, \quad i = \sqrt{-1}. \quad (10)$$

**Assumption 1**  $\Phi(\omega)$  is a rational power spectral density with the order of  $2r$  with  $r \in \mathbb{N}$ .

Under Assumption 1, the spectral factorization technique, e.g., [2, 15], can be applied to (10) and there exists a rational transfer function  $W$  such that

$$\Phi(\omega) = W(e^{i\omega})W(e^{-i\omega}). \quad (11)$$

From the realization theory of linear systems e.g., [8] and the transfer function  $W(e^{i\omega})$ , for each location  $p_i$ , the corresponding discrete-time state-space model realization of a zero mean Gaussian process with the covariance function (9) can be derived by

$$\begin{aligned} s_{i,j} &= F_D s_{i,j-1} + G_D w_{i,j-1}, \quad s_{i,0} \sim \mathcal{N}(0, \Sigma_0), \\ z_{i,j} &= H_D s_{i,j}, \quad j = 1, \dots, \end{aligned} \quad (12)$$

with  $i, i' = 1, \dots, M, j, j' = 1, \dots, N + N_T$ ,

$$\mathbb{E}[z_{i,j'} z_{i,j}] = k_t(\tau; \alpha_t), \quad \mathbb{E}[z_{i,j} z_{i',j}] = 0, \quad (13)$$

where  $F_D \in \mathbb{R}^{r \times r}, G_D \in \mathbb{R}^r$  and  $H_D \in \mathbb{R}^{1 \times r}$  are the system matrix, the input matrix and the output matrix, respectively,  $s_{i,j} \in \mathbb{R}^r$  is the state vector of the  $i$ th location at the  $j$ th time instant with  $s_{i,0}$  and  $s_{i',0}$  ( $i \neq i'$ ) being independent from each other,  $w_{i,j} \in \mathbb{R}$  is white Gaussian noise with zero mean and unit variance,  $\Sigma_0$  is the solution of the discrete-time Lyapunov equation  $\Sigma_0 = F_D \Sigma_0 F_D^T + G_D G_D^T \in \mathbb{R}^{r \times r}$  and  $\mathbb{E}(\cdot)$  is the mathematical expectation.

Then we define that

$$\chi_j = [f(p_1, t_j), \dots, f(p_M, t_j)]^T \in \mathbb{R}^M, \quad (14)$$

and according to Assumption 1 and (4), its covariance matrix is

$$\mathbb{E}[\chi_{j'} \chi_j^T] = K_s(\alpha_s) k_t(\tau; \alpha_t) \in \mathbb{R}^{M \times M}, \quad (15)$$

where  $K_s(\alpha_s)$  is defined in (8c). We let

$$z_j = [z_{1,j}, \dots, z_{M,j}]^T \in \mathbb{R}^M, \quad (16)$$

and then with (13) and (15), we obtain

$$\chi_j = K_s(\alpha_s)^{1/2} z_j, \quad (17)$$

where  $K_s(\alpha_s)^{1/2}$  is the ‘‘square root’’ of  $K_s(\alpha_s)$  defined in (20). With (15)-(17), we rewrite (12) as follows

$$s_j = F s_{j-1} + G w_{j-1}, \quad s_0 \sim \mathcal{N}(0, I_M \otimes \Sigma_0), \quad (18a)$$

$$\chi_j = H s_j, \quad j = 1, 2, \dots, \quad (18b)$$

where  $s_j = [s_{1,j}^T, \dots, s_{M,j}^T]^T \in \mathbb{R}^{Mr}$ ,  $F = I_M \otimes F_D \in \mathbb{R}^{Mr \times Mr}$ ,  $G = I_M \otimes G_D \in \mathbb{R}^{Mr \times M}$ ,  $H = K_s(\alpha_s)^{1/2} (I_M \otimes H_D) \in \mathbb{R}^{M \times Mr}$  and  $w_j = [w_{1,j}, \dots, w_{M,j}]^T \in \mathbb{R}^M$ .

According to (17) and (15), the state-space model (18) is a realization of the Gaussian process (3). Then the model (1) can be accordingly rewritten as follows

$$s_j = F s_{j-1} + G w_{j-1}, \quad s_0 \sim \mathcal{N}(0, I_M \otimes \Sigma_0), \quad (19a)$$

$$y_j = H s_j + v_j, \quad j = 1, 2, \dots, \quad (19b)$$

where  $v_j = [v_{1,j}, \dots, v_{M,j}]^T \in \mathbb{R}^M$ , with  $v_j \sim \mathcal{N}(0, \sigma^2 I_M)$ ,  $y_j$  is defined in (8d), and  $w_j$  and  $v_j$  are independent for any  $j = 1, \dots, N + N_T$ .

Then the spatial-temporal function estimation and prediction problem can be converted to a Kalman filtering, smoothing and prediction problem for (19) and the corresponding implementations has computational complexity  $\mathcal{O}(NM^3)$ , same as the ones in e.g., [7, 26, 27].

**Remark 1** Note that Section 3.1 and [27, Proposition 2] use two different routes to derive the discrete-time state-space model realization (19) of the spatial-temporal Gaussian process. It is not hard to show that they are equivalent in theory, but they are different in implementation. In particular, the discretization technique used in [27] includes solving an integral involving the matrix exponential, which needs to be handled carefully and if otherwise, numerical problem may occur, see e.g., [28] and the references therein. Therefore, the route in Section 3.1 is preferable in practice, because no discretization of continuous-time state-space model is involved and thus possible numerical problems are avoided. Moreover, in practice one can design directly the discrete-time simulation-induced kernel [9] based on the prior knowledge, which is represented in a state-space model form.

### 3.2 A Transformed State-space Model Realization

In order to further reduce the computational complexity in terms of  $M$ , it is useful to explore the Kronecker

structure of the system, input and output matrices of the state-space model (19) and perform a coordinate and an output transformation to (19).

Firstly, we denote the singular value decomposition (SVD) of the spatial kernel matrix  $K_s(\alpha_s)$  and its ‘‘square root’’  $K_s(\alpha_s)^{\frac{1}{2}}$  as follows

$$K_s(\alpha_s) = \Lambda D \Lambda^T, K_s(\alpha_s)^{\frac{1}{2}} = \Lambda D^{\frac{1}{2}} \Lambda^T, \quad (20)$$

where  $D \in \mathbb{R}^{M \times M}$  is a diagonal matrix and its main diagonals are singular values of  $K_s(\alpha_s)$ ,  $D^{\frac{1}{2}}$  is a diagonal matrix with the square root of diagonals of  $D$  and  $\Lambda$  is an orthogonal matrix, i.e.  $\Lambda \Lambda^T = \Lambda^T \Lambda = I_M$ .

Then for  $j = 0, 1, \dots$ , we introduce a state transform

$$\begin{aligned} x_j &= (\Lambda^T \otimes I_r) s_j \iff s_j = (\Lambda \otimes I_r) x_j, \\ x_{j+1} &= (\Lambda^T \otimes I_r) (F s_j + G w_j) \\ &= (\Lambda^T \otimes I_r) F (\Lambda \otimes I_r) x_j + G \Lambda^T w_j \\ &= (I_M \otimes F_D) x_j + G \Lambda^T w_j, \end{aligned} \quad (21)$$

where  $(\Lambda^T \otimes I_r) G = (\Lambda^T \otimes G_D) = (I_M \otimes G_D) (\Lambda^T \otimes 1) = G \Lambda^T$ , and an output transform for  $j = 1, 2, \dots$ ,

$$\begin{aligned} l_j &= \Lambda^T y_j \\ &= \Lambda^T (H s_j + v_j) \\ &= \Lambda^T H (\Lambda \otimes I_r) x_j + \Lambda^T v_j \\ &= \Lambda^T K_s(\alpha_s)^{\frac{1}{2}} (I_M \otimes H_D) (\Lambda \otimes I_r) x_j + \Lambda^T v_j \\ &= D^{\frac{1}{2}} \Lambda^T (I_M \otimes H_D) (\Lambda \otimes I_r) x_j + \Lambda^T v_j \\ &= D^{\frac{1}{2}} (I_M \otimes H_D) x_j + \Lambda^T v_j, \end{aligned} \quad (22)$$

where the last equation is true because  $\Lambda^T (I_M \otimes H_D) (\Lambda \otimes I_r) = (\Lambda^T \otimes 1) (\Lambda \otimes H_D) = I_M \otimes H_D$ . Then the state-space model (19) is transformed to

$$x_j = \bar{F} x_{j-1} + G \bar{w}_{j-1}, x_0 \sim \mathcal{N}(0, I_M \otimes \Sigma_0), \quad (23a)$$

$$l_j = \bar{H} x_j + \bar{v}_j, j = 1, 2, \dots, \quad (23b)$$

where  $\bar{F} \in \mathbb{R}^{M_r \times M_r}$ ,  $\bar{H} \in \mathbb{R}^{M \times M_r}$  and the covariance of  $x_0$  are computed as

$$\bar{F} = I_M \otimes F_D, \bar{H} = D^{\frac{1}{2}} (I_M \otimes H_D), \quad (24a)$$

$$I_M \otimes \Sigma_0 = (\Lambda \otimes I_r) (I_M \otimes \Sigma_0) (\Lambda^T \otimes I_r), \quad (24b)$$

$\bar{w}_j = \Lambda^T w_j \sim \mathcal{N}(0, I_M)$  and  $\bar{v}_j = \Lambda^T v_j \sim \mathcal{N}(0, \sigma^2 I_M)$ . We denote the transformed output vector and its covariance matrix by  $L$  and  $\bar{\Sigma}(\alpha)$ , respectively, which are described by

$$L = [l_1^T, \dots, l_N^T] = (I_N \otimes \Lambda^T) Y, \quad (25a)$$

$$\bar{\Sigma}(\alpha) = \text{COV}[L, L] = (I_N \otimes \Lambda^T) \Sigma(\alpha) (I_N \otimes \Lambda). \quad (25b)$$

### 3.3 Kalman Filter Based Estimation and Prediction

Firstly, we define the estimate  $\hat{x}_{j|m}$  and its covariance matrix  $\bar{\Sigma}_{j|m}$  for  $j = 1, 2, \dots, m = 0, 1, \dots, N$  as

$$\hat{x}_{j|m} = \mathbb{E}[x_j | l_{0:m}], \quad (26a)$$

$$\bar{\Sigma}_{j|m} = \mathbb{E}[(x_j - \hat{x}_{j|m})(x_j - \hat{x}_{j|m})^T | l_{0:m}], \quad (26b)$$

where  $l_0$  is a null vector and

$$l_{0:m} = \{l_0, \dots, l_m\}. \quad (27)$$

Then the Kalman filter for (23) can be expressed as

$$\bar{e}_j = l_j - \bar{H} \hat{x}_{j-1}, \quad (28a)$$

$$\bar{E}_j = \text{COV}[\bar{e}_j, \bar{e}_j] \quad (28b)$$

$$= \bar{H} \bar{\Sigma}_{j-1} \bar{H}^T + \sigma^2 I_M, \quad (28c)$$

$$\hat{x}_{j|j} = \hat{x}_{j-1|j-1} + \bar{\Sigma}_{j-1} \bar{H}^T \bar{E}_j^{-1} \bar{e}_j, \quad (28d)$$

$$\bar{\Sigma}_{j|j} = \bar{\Sigma}_{j-1|j-1} - \bar{\Sigma}_{j-1} \bar{H}^T \bar{E}_j^{-1} \bar{H} \bar{\Sigma}_{j-1}, \quad (28e)$$

$$\hat{x}_{j+1|j} = \bar{F} \hat{x}_{j|j}, \quad (28f)$$

$$\bar{\Sigma}_{j+1|j} = \bar{F} \bar{\Sigma}_{j|j} \bar{F}^T + Q, \quad (28g)$$

where  $\bar{e}_j$  is known as the innovation,

$$Q = I_M \otimes (G_D G_D^T) \quad (29)$$

and the iterative algorithm starts from  $j = 1$ .

For the purpose of function estimation, we apply Kalman smoother as follows

$$\bar{\Sigma}_{j|j} = \bar{\Sigma}_{j-1|j-1} - \bar{\Sigma}_{j-1} \bar{H}^T \bar{E}_j^{-1} \bar{H} \bar{\Sigma}_{j-1}, \quad (30a)$$

$$\bar{J}_j = \bar{\Sigma}_{j|j} \bar{F}^T \bar{\Sigma}_{j+1|j}^{-1}, \quad (30b)$$

$$\hat{x}_{j|N} = \hat{x}_{j|j} + \bar{J}_j (\hat{x}_{j+1|N} - \bar{F} \hat{x}_{j|j}), \quad (30c)$$

$$\bar{\Sigma}_{j|N} = \bar{\Sigma}_{j|j} + \bar{J}_j (\bar{\Sigma}_{j+1|N} - \bar{\Sigma}_{j+1|j}) \bar{J}_j^T, \quad (30d)$$

$$\hat{f}_{j|N} = \Lambda \bar{H} \hat{x}_{j|N}, j = N - 1, \dots, 1, \quad (30e)$$

where  $\hat{f}_{j|N} = \mathbb{E}[\chi_j | l_{0:N}]$  with  $\chi_j$  defined in (14).

For the purpose of function prediction, we apply the Kalman predictor as follows

$$\hat{x}_{j|N} = \bar{F} \hat{x}_{j-1|N}, \quad (31a)$$

$$\bar{\Sigma}_{j|N} = \bar{F} \bar{\Sigma}_{j-1|N} \bar{F}^T + Q, \quad (31b)$$

$$\hat{f}_{j|N} = \Lambda \bar{H} \hat{x}_{j|N}, j = N + 1, \dots, N + N_T, \quad (31c)$$

where  $\hat{f}_{j|N} = \mathbb{E}[\chi_j | l_{0:N}]$  is the prediction of  $\chi_j$  at  $j$ th time instant.

### 3.4 Hyper-parameter Estimation

Based on the Kalman filter (28), it is possible to propose efficient implementation algorithms for the MLM, GCV and SURE methods.

**Lemma 1** [6, p. 302, Properties of the Innovation Sequence] For  $j = 1, \dots, N$ , the innovation  $\bar{e}_j$  in (28a) can be represented as a linear function of  $l_{0:j}$  in (27), i.e.

$$\bar{e}_1 = l_1, \quad (32a)$$

$$\bar{e}_j = l_j - \sum_{i=1}^{j-1} b_{j,i} l_i, \text{ for } j = 2, \dots, N, \quad (32b)$$

$$\text{COV}[\bar{e}_j, \bar{e}_{j'}] = 0, \text{ for } j' = 1, \dots, N \text{ and } j \neq j', \quad (32c)$$

where  $b_{j,i} \in \mathbb{R}$  is the corresponding coefficient for  $i = 1, \dots, j-1$ .

**Proposition 1** Let

$$\Theta = [\bar{e}_1^T, \dots, \bar{e}_N^T]^T, \quad (33a)$$

$$\Psi = \text{COV}[\Theta, \Theta], \quad (33b)$$

where  $\bar{e}_j$ ,  $j = 1, \dots, N$ , are defined in (28a). Then following Lemma 1,  $\Theta$  and  $\Psi$  can be rewritten as

$$\Theta = \Gamma L, \quad (34a)$$

$$\Psi = \Gamma \bar{\Sigma}(\alpha) \Gamma^T, \quad (34b)$$

$$= \text{blkdiag}(\bar{E}_1, \dots, \bar{E}_N), \quad (34c)$$

where  $L$  is defined in (25a),  $\Gamma \in \mathbb{R}^{NM \times NM}$  is a lower uni-triangular matrix with  $|\Gamma| = 1$  and for  $j, i = 1, \dots, NM$ , the  $(j, i)$ th element of  $\Gamma$  is

$$[\Gamma]_{ji} = \begin{cases} 0, & j < i, \\ 1, & j = i, \\ b_{j,i}, & j > i, \end{cases} \quad (35)$$

and  $\text{blkdiag}(\bar{E}_1, \dots, \bar{E}_N)$  is a block diagonal matrix with  $\bar{E}_1, \dots, \bar{E}_N$ , defined in (28b), on the main diagonals.

By Proposition 1, the cost function of the MLM method (5) can be calculated as shown in the proposition below.

**Proposition 2** The cost function of the MLM method (5) can be computed by using

$$\log |\Sigma(\alpha)| = \sum_{j=1}^N \log |\bar{E}_j|, \quad (36a)$$

$$Y^T \Sigma^{-1}(\alpha) Y = \sum_{j=1}^N \bar{e}_j^T \bar{E}_j^{-1} \bar{e}_j. \quad (36b)$$

Then by Propositions 1 and 2, the cost functions of the GCV and SURE methods can be calculated as shown in the following proposition.

**Proposition 3** The cost functions of the GCV method (6) and the SURE method (7) can be computed by using

$$S = \sigma^4 \sum_{j=1}^N [\bar{e}_j^T \bar{E}_j^{-1} (\bar{H} \bar{P}_{j|j-1} \bar{H}^T + I_M) \bar{E}_j^{-1} \bar{e}_j + 2 \bar{\zeta}_{j|j-1}^T \bar{H} \bar{E}_j^{-1} \bar{e}_j], \quad (37)$$

$$\delta = MN - \sigma^2 \sum_{j=1}^N \text{trace} [\bar{E}_j^{-1} (\bar{H} \bar{P}_{j|j-1} \bar{H}^T + I_M)],$$

where  $\bar{\zeta}_{j|j-1}$  and  $\bar{P}_{j|j-1}$  can be computed recursively:

- for  $j = 1$ ,  $\bar{\zeta}_{1|0} = 0 \in \mathbb{R}^{Mr}$  and  $\bar{P}_{1|0} = 0 \in \mathbb{R}^{Mr \times Mr}$ ;
- for  $j = 2, \dots, N$ ,

$$\begin{aligned} \bar{\zeta}_{j|j-1} &= \bar{F} \bar{\zeta}_{j-1|j-2} + \bar{F} \bar{P}_{j-1|j-2} \bar{H}^T \bar{E}_j^{-1} \bar{e}_j \\ &\quad - \bar{F} \bar{\Sigma}_{j-1|j-2} \bar{H}^T \bar{E}_j^{-1} (\bar{H} \bar{P}_{j|j-1} \bar{H}^T + I_M) \bar{E}_j^{-1} \bar{e}_j \\ &\quad - \bar{F} \bar{\Sigma}_{j-1|j-2} \bar{H}^T \bar{E}_j^{-1} \bar{H} \bar{\zeta}_{j-1|j-2}, \end{aligned} \quad (38)$$

$$\begin{aligned} \bar{P}_{j|j-1} &= \bar{F} \bar{P}_{j-1|j-2} \bar{F}^T - \bar{F} \bar{P}_{j-1|j-2} \bar{H}^T \bar{E}_j^{-1} \bar{H} \bar{\Sigma}_{j-1|j-2} \bar{F}^T \\ &\quad - \bar{F} \bar{\Sigma}_{j-1|j-2} \bar{H}^T \bar{E}_j^{-1} \bar{H} \bar{P}_{j-1|j-2} \bar{F}^T \\ &\quad - \bar{F} \bar{\Sigma}_{j-1|j-2} \bar{H}^T \bar{E}_j^{-1} (\bar{H} \bar{P}_{j-1|j-2} \bar{H}^T + I_M) \\ &\quad \bar{E}_j^{-1} \bar{H} \bar{\Sigma}_{j-1|j-2} \bar{F}^T. \end{aligned} \quad (39)$$

### 3.5 Summary of the Implementation Algorithm and Its Computational Complexity Analysis

The proposed implementation, as shown in Sections 3.2-3.4, can be summarized in Algorithm 1 below. To analyze the computational complexity of Algorithm 1, it should be noted that the dimension  $r$  of the state-space model (12) is determined by the temporal kernel (9) and is irrespective of, and often much smaller than,  $M$  and  $N$ , and thus in what follows, we ignore  $r$  and moreover, let  $N_T = N$  in the analysis for brevity.

**Theorem 1** The proposed implementation, as shown in Algorithm 1, has computational complexity  $\mathcal{O}(M^3 + NM^2)$ . In particular,

- the state-space model transformation (20), (24a) and (25a) has computational complexity  $\mathcal{O}(M^3 + NM^2)$ ;
- the Kalman filter (28) has computational complexity  $\mathcal{O}(NM)$ ;
- the evaluation of the cost functions of the MLM method (5), the GCV method (6) and the SURE method (7) has computational complexity  $\mathcal{O}(NM)$ ;
- the Kalman smoother (30) and Kalman predictor (31) have computational complexity  $\mathcal{O}(NM^2)$ .

---

**Algorithm 1** The Proposed Implementation
 

---

**Input:** data  $\{p_i, t_j, y_{i,j}\}_{i=1, j=1}^{M, N+N_T}$ , kernels  $k_t(t_j, t_j'; \alpha_t)$ ,  $k_s(p_i, p_i'; \alpha_s)$

**Output:**  $\hat{f}_{j|N}$  and  $\bar{\Sigma}_{j|N}$  for  $j = 1, \dots, N + N_T$ .

**Step 1:** State-space model derivation  
Derive (19);

**Step 2:** State-space model transformation  
Calculate (20), (24a) and (25a);

**Step 3:** Hyper-parameter Estimation

- Kalman filter  
Calculate (28);

**if** use the MLM method (5) **then**  
Calculate (36);  
**end**

**if** use the GCV method (6) **then**  
Calculate (38), (39), (A.8a), (A.8b) and (A.6);  
**end**

**if** use the SURE method (7) **then**  
Calculate (38), (39), (A.8a), (A.8b) and (A.6);  
**end**

**Step 4:** Function estimation and prediction

- Kalman smoother for estimation  
Calculate (30);
- Kalman predictor for prediction  
Calculate (31);

---

**Remark 2** For spatial-temporal data with large  $N$  and moderately large  $M$ , to reduce the computational complexity of the MLM method (5), the GCV method (6) and the SURE method (7), it is suggested to use derivative-free optimization algorithms or algorithms that only require numerical gradient, approximated by finite difference of the cost function of the optimization problems involved. With such optimization algorithms, solving the MLM method (5), the GCV method (6) and the SURE method (7) only involves the state-space model transformation, the Kalman filter and the evaluation of the cost function of the optimization problems and thus has computational complexity  $\mathcal{O}(M^3 + NM^2)$ .

## 4 Applications

In this section, we illustrate the proposed implementation over applications in weather data prediction and spatially-distributed system identification.

### 4.1 Computing Platform

Firstly, we introduce our computing platform in Fig. 1, which consists of 1 server and 2 GPUs:

- Server 1: Intel(R) Xeon(R) Platinum 8168 2.7GHz CPU×2 (48 cores), 64GB×24=1.48TB RAM,
- GPU: NVIDIA V100 ×2, 16GB RAM.

It is worth to note that many computations in the proposed implementation can be parallelized. For example,

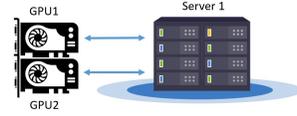


Fig. 1. Computing platform

the creation of the spatial kernel matrix  $K_s(\alpha_s)$ , the SVD of  $K_s(\alpha_s)$ , the output transformation in (22), the computation of (30e) and (31c). Then by using the parallel computing structure of the computing platform and the parallel computing toolbox in MATLAB, the proposed implementation can be made more efficient.

### 4.2 Weather Data Prediction

#### 4.2.1 Weather Data Sets

We consider the following two weather data sets.

- 1) **Colorado Precipitation Data:** This data set has been tested in e.g., [17, 27], contains monthly precipitation data between 1895 and 1997 from 367 weather stations in Colorado, USA<sup>1</sup>. The data set contains in total 1236 time instants and 367 locations (stations) located in a rectangular longitude/latitude region  $[109.5^\circ\text{W}, 101^\circ\text{W}] \times [36.5^\circ\text{N}, 41.5^\circ\text{N}]$ . We treat the data in 1895 – 1995 as the training data, and the data in 1996 – 1997 as the test data, that is, we have  $t_j = jT_s$ ,  $j = 1, \dots, N + N_T$  with  $T_s = 1$  month,  $N = 1212$ ,  $N_T = 24$  and  $M = 367$ . This data set contains in total 453,612 data points.
- 2) **GHCN Temperature Data:** This data set is obtained from the Global Historical Climatology Network (GHCN), and contains daily average temperatures collected from over ten thousands weather stations over the world [21]. We first choose 4000 stations with most complete data records from the 301th day of 1999 to the 300th day of 2018. Then we take out those stations with daily average temperature over  $80^\circ\text{C}$  or under  $-80^\circ\text{C}$  and there are 3955 locations left. The data set contains in total 6940 times instants and 3955 locations (stations) and we treat the data in the former 18 years as the training data, and the data in the last year as the test data, that is, we have  $t_j = jT_s$ ,  $j = 1, \dots, N + N_T$  with  $T_s = 1$  day,  $N = 6575$ ,  $N_T = 365$  and  $M = 3955$ . This data consist of more than 27 million data points and is much larger than the Colorado precipitation data.

#### 4.2.2 Kernel Design

In this section, we design a kernel for both the Colorado precipitation data and GHCN temperature data.

<sup>1</sup> <https://www.image.ucar.edu/Data/US.monthly.met/CO.shtml>.

The kernel design problem here is tricky, because the designed kernel should on the one hand incorporate the prior knowledge on the underlying spatial-temporal function to be estimated and on the other hand has state-space model realizations.

We first consider the spatial kernel design. Since the precipitation and the temperature are diffusion processes, the spatial prior knowledge is that for two locations, the closer the two locations, the larger the correlation between their weather data, and thus the squared exponential (SE) kernel is often adopted, e.g., [27, 29],

$$k_{\text{SE}}(p_i, p_{i'}; \alpha_{se}) = \exp\left(-\|p_i - p_{i'}\|_2^2 / \alpha_{se}\right), \quad (40)$$

where  $\alpha_{se} > 0$ , and for the Colorado precipitation data,  $p_i \in \mathbb{R}^2$  and its components are the longitude and latitude of the location, respectively, and for the GHCN temperature data,  $p_i \in \mathbb{R}^3$  and its components are the earth-centered earth-fixed (ECEF) coordinates of the locations considered and the units are in 10 kilometers.

Then we consider the temporal kernel design. To capture the periodicity of the weather data, an intuitive way is to use the periodic kernel, e.g., [29]

$$k_{\text{per}}(\tau; \delta_t, c_t) = \delta_t \exp\left\{-2c_t [\sin(\pi \mathbf{f} \tau)]^2\right\}, \quad (41)$$

where  $c_t > 0$ ,  $\delta_t > 0$  are the hyper-parameters, and  $\mathbf{f} \in \mathbb{R}$  is the period of the weather data. However, the periodic kernel (41) does not have a proper PSD and thus has no state-space model realization. To overcome this difficulty, we first consider the Taylor expansion of  $\exp(x)$  at  $x = 0$  to the second-order, then replace  $x$  by  $-2c_t[\sin(\pi \mathbf{f} \tau)]^2$  in the expansion, and finally, multiply it by an exponential kernel  $k_{\text{EXP}}(\tau) = e^{-|\tau|/\sigma_t}$  and obtain the following positive definite kernel

$$\begin{aligned} & k_{\text{TE2}}(\tau; \delta_t, c_t) k_{\text{EXP}}(\tau) \\ &= \delta_t \left[ \left(1 - c_t + \frac{3}{4}c_t^2\right) + (c_t - c_t^2) \cos(2\pi \mathbf{f} |\tau|) \right. \\ & \quad \left. + \frac{c_t^2}{4} \cos(4\pi \mathbf{f} |\tau|) \right] \exp\left(-\frac{|\tau|}{\sigma_t}\right), \end{aligned} \quad (42)$$

where  $\delta_t > 0$  and  $c_t \in (0, 1)$  are the hyper-parameters,  $c_t \in (0, 1)$  is imposed to guarantee that (42) is positive semidefinite. The derivation of the state-space model of (42) is included in Appendix A.5. Here, it should be noted that both  $\mathbf{f}$  and  $\sigma_t$  are not hyper-parameters:  $\mathbf{f}$  is chosen to be  $\mathbf{f} = 1/12$  for the Colorado precipitation data and  $\mathbf{f} = 1/365.3$  for the GHCN temperature data due to the periodicity of the data, and  $\sigma_t$  is chosen to be  $\sigma_t = 5 \times 10^3$  such that the exponential kernel  $k_{\text{EXP}}(\tau) = e^{-|\tau|/\sigma_t}$  has a negligible effect. Moreover, to describe the slight temporal variation of the data, we further include

a Matérn kernel, e.g., [29], i.e.,

$$k_{\text{Matern}}(\tau; h_t, \theta_t) = h_t \left(1 + \frac{\sqrt{3}|\tau|}{\theta_t}\right) \exp\left(-\frac{\sqrt{3}|\tau|}{\theta_t}\right), \quad (43)$$

where  $h_t, \theta_t > 0$  are the hyper-parameters of (43). Then we can obtain the following temporal kernel

$$k_t(\tau; \alpha_{tm}) = k_{\text{TE2}}(\tau; \delta_t, c_t) k_{\text{EXP}}(\tau) + k_{\text{Matern}}(\tau; h_t, \theta_t), \quad (44)$$

where  $\alpha_{tm} = [\delta_t, c_t, h_t, \theta_t]^T$  with

$$0.01\delta_t \leq h_t \leq 0.1\delta_t, \quad (45)$$

which is enforced to guarantee that the Matérn kernel (43) describes the slight temporal variation of the data.

**Remark 3** *Beside the spatial prior knowledge considered above, it is interesting to note that the spatial prior knowledge considered in [32] is that the edges in the graphical model are sparse, where the graphical model is due to the existence of a number of modules with a graphical structure, and that each module has a number of nodes sharing the same graphical structure, and thus, a sparsity inducing kernel/regularization was designed accordingly. It is also interesting to mention that the following kernel*

$$k_{\text{PD}}(\tau; \delta_t, \sigma_t) = \delta_t \cos(2\pi \mathbf{f} |\tau|) e^{-\frac{|\tau|}{\sigma_t}}, \quad (46)$$

where  $\delta_t, \sigma_t > 0$  are hyper-parameters and  $\mathbf{f} = 1/12$ , is chosen in [17, 27] as the temporal kernel.

#### 4.2.3 Hyper-parameter Estimation and Function Prediction

For the two data sets and designed kernels, we use the MLM, GCV and SURE methods, as shown in Section 3.4, to estimate the hyper-parameter  $\alpha = [\alpha_t^T, \alpha_s^T]^T$ . Moreover, for the MLM method, the noise variance  $\sigma^2$  is treated as an additional hyper-parameter, i.e.,  $\alpha = [\alpha_t^T, \alpha_s^T, \sigma^2]^T$ , and its estimate is then used for the SURE and GCV methods. With the estimated hyper-parameter, we can further run the Kalman filter, smoother and predictor in Section 3.3 to compute the function prediction  $\hat{f}_T = [\hat{f}_{N+1|N}, \dots, \hat{f}_{N+N_T|N}]$ .

The function `fmincon` in Matlab, using the interior-point algorithm with numerical gradient approximated by finite difference of the cost function, is applied to solve (5), (6) or (7). Since the selection of initial points is significant for the search of “good” local minima, the following way is used to find a “good” local minimum:

- 1) for each component of the hyper-parameter, we select a set of initial points and thus obtain a grid of initial points of the hyper-parameter;

- 2) calculate the cost functions over the grid of initial points;
- 3) select 5 initial points corresponding to the smallest 5 values of the cost function;
- 4) use the function `fmincon` with selected 5 initial points to solve the optimization problem involved in the hyper-parameter estimation, respectively;
- 5) choose the optimal solution with the smallest value of the cost function as the optimal hyper-parameter estimate.

To assess how good the prediction  $\hat{f}_T$  is at  $jT_s$  for  $j = N + 1, \dots, N + N_T$ , we use the measure of fit, e.g., [20],

$$\text{fit}_j = 100 \left( 1 - \frac{\|\hat{f}_{j|N} - y_j\|_2}{\|y_j - \bar{y}_j\|_2} \right), \bar{y}_j = \frac{1}{M} \sum_{i=1}^M y_{i,j}. \quad (47)$$

The maximum of  $\text{fit}_j$  is 100, meaning a perfect match between  $\hat{f}_{j|N}$  and  $y_j$ . The average prediction fit over the test data set is defined as

$$\bar{\text{fit}} = \frac{1}{N_T} \sum_{j=N+1}^{N+N_T} \text{fit}_j. \quad (48)$$

#### 4.2.4 Filling the Missing Data

The Colorado precipitation and GHCN temperature data contain 58.39% and 3.3% missing data, respectively, and we need to fill the missing data before running simulations. To this goal, we first split the spatial-temporal data into  $M$  temporal data sets  $\{p_1, t_j, y_{1,j}\}_{j=1}^{N+N_T}, \dots, \{p_M, t_j, y_{M,j}\}_{j=1}^{N+N_T}$  according to the  $M$  locations. For each temporal data set, the temporal kernel (46) or (44) is applied, respectively. Then for each  $i = 1, \dots, M$ , we use the MLM method to estimate the corresponding hyper-parameter and in particular, if  $y_{i,j}$  is missing for some  $j = 1, \dots, N + N_T$ , then no measurement update is needed, i.e., (28d) and (28e) should be replaced by

$$\hat{x}_{i|j} = \hat{x}_{i|j-1}, \quad \bar{\Sigma}_{i|j} = \bar{\Sigma}_{i|j-1}, \quad (49)$$

respectively, e.g., [1]. Finally, with the obtained hyper-parameter estimate, the Kalman smoother (30) is used to fill the missing data.

**Remark 4** *The above treatment of the missing data implicitly assumes that for  $i, i' = 1, \dots, M$  and  $j, j' = 1, \dots, N + N_T$ , if  $i \neq i'$ ,  $y_{i,j}$  and  $y_{i',j'}$  are independent. The treatment in [27] does not rely on this assumption and thus is more general but with the price of higher computational complexity.*

#### 4.2.5 Illustration of Computational Efficiency

Firstly, we consider the Colorado precipitation data and choose (40) as the spatial kernel and (46) as the temporal

kernel, and then we evaluate the cost functions of the MLM method (5), GCV method (6), and SURE method (7) for 10 times. The average computing time of the cost functions for the proposed implementation and the one in [17, 27] are shown in the Table 1, which shows that, our proposed implementation is over 300 and 200 times faster than the one in [17, 27] for the MLM method, and GCV and SURE methods, respectively.

Table 1

The average computing time (in second) of the cost functions of the MLM method (5), GCV method (6) and SURE method (7) for the Colorado precipitation data.

Implementation	Proposed	in [17, 27]
MLM method	0.6485	197.3668
GCV method	1.5983	332.9270
SURE method	1.5904	331.9392

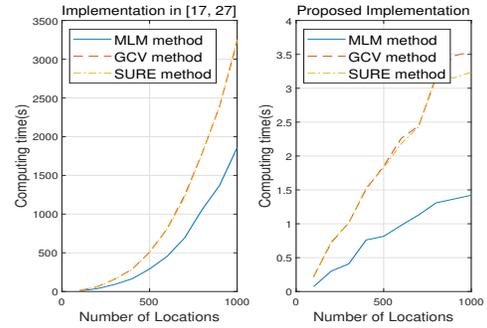


Fig. 2. The average computation time (in second) of the cost functions of the MLM method (5), GCV method (6) and SURE method (7) for the GHCN temperature data with  $N = 800$  and  $M = 100, 200, \dots, 1000$ , respectively, where our proposed implementation and the one in [17, 27] are shown on the right and left panels, respectively.

Secondly, we consider the GHCN temperature data but only use part of it, because the implementation in [17, 27] is too expensive to be applied to the full data. In particular, we only use the first 800 time instants and 1000 locations, i.e.,  $\{p_i, t_j, y_{i,j}\}_{i=1, j=1}^{1000, 800}$ . Then we choose (40) and (44) as the spatial kernel and temporal kernel, respectively, and evaluate the cost functions of the MLM method (5), GCV method (6), and SURE method (7) for 10 times with  $N = 800$  and  $M = 100, \dots, 1000$ , respectively. The average computing time of the cost functions for the proposed implementation and the one in [17, 27] are shown Fig. 2, which shows that, our proposed implementations is more efficient than the one in [17, 27], as the number of the locations increases. It is worth to mention that for the full GHCN temperature data, our proposed implementation has the average computing time 30.2, 68.2 and 67.8 seconds, for the cost functions of the MLM, GCV, and SURE methods, respectively.

#### 4.2.6 Illustration of Prediction Performance

For the Colorado precipitation data, the prediction fits (47), the average prediction fits (48), and the opti-

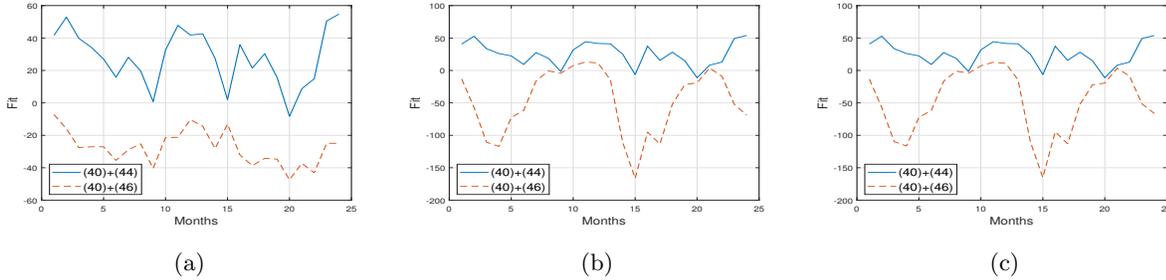


Fig. 3. Profile: Monthly prediction fits (47) for the Colorado precipitation data using two kernel combinations: (40)+(44) and (40)+(46), and three hyper-parameter estimation methods, respectively. Panel (a): Monthly prediction fits using the MLM method. Panel (b): Monthly prediction fits using the GCV method. Panel (c): Monthly prediction fits using the SURE method.

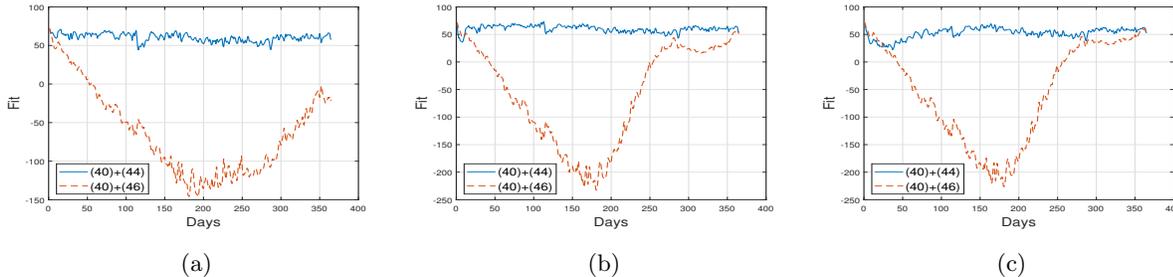


Fig. 4. Profile: Daily prediction fits (47) for the GHCN temperature data using two kernel combinations: (40)+(44) and (40)+(46), and three hyper-parameter estimation methods, respectively. Panels (a): Daily prediction fits using the MLM method. Panel (b): Daily prediction fits using the GCV method. Panel (c): Daily prediction fits using the SURE method.

mal hyper-parameters using two kernel combinations: (40)+(44) and (40)+(46), and three hyper-parameter estimation methods: MLM, GCV and SURE, are shown in Fig. 3, Tables 2 and 3, respectively. It is worth to stress that (40)+(46) was used in [17, 27].

Fig. 3 and Table 2 show that for the same hyper-parameter estimation method, the prediction fits (47) and the average prediction fits (48) obtained by the kernel combination (40)+(44) are all larger than those by (40)+(46), indicating that the temporal kernel (44) can better describe the Colorado precipitation data than (46) used in [17, 27]. Table 2 also shows that the kernel combination (40)+(44) with hyper-parameters estimated by the MLM method gives the best average prediction fit 28.25. One may wonder why this fit is not so good and the reason is perhaps due to that 58.39% of the Colorado precipitation data are missing.

Table 2

The average prediction fits (48) for the Colorado precipitation data using two kernel combinations, where the hyper-parameters are estimated by MLM, GCV and SURE methods, respectively. The values in parentheses are the corresponding smallest prediction fits.

Kernel	(40)+(44)	(40)+(46)
MLM	28.25 (-8.31)	-27.54 (-47.39)
GCV	25.66 (-11.24)	-47.56 (-166.35)
SURE	25.66 (-11.24)	-47.38 (-165.38)

For the GHCN temperature data, the prediction fits (47), the average prediction fits (48), and the optimal hyper-parameters using two kernel combinations: (40)+(44) and (40)+(46), and three hyper-parameter estimation methods: MLM, GCV and SURE, are shown in Fig. 4, Tables 4 and 5, respectively. It is worth to stress that (40)+(46) was used in [17, 27].

Fig. 4 and Table 4 show that for the same hyper-parameter estimation method, the prediction fits (47) and the average prediction fits (48) obtained by the kernel combination (40)+(44) are most of time larger than those by (40)+(46), indicating that the temporal kernel (44) can better describe the GHCN temperature data than (46) used in [17, 27]. Table 4 also shows, among three hyper-parameter estimation methods, the MLM method gives the best average prediction fit 60.34 and the corresponding smallest prediction fit 44.12. Moreover, Fig. 4 also shows that the prediction fit of the GCV and SURE methods drop down quickly in the beginning and then go up again, while the MLM method can avoid such drop.

### 4.3 Spatially-distributed System Identification

In this section, we consider the identification of *spatially-distributed* system, e.g. [18], which is a class of distributed parameter systems.

Table 3

The hyper-parameter and the corresponding optimal values of the cost functions for the Colorado precipitation data using two kernel combinations: (40)+(44) and (40)+(46), and three hyper-parameter estimation methods, respectively.

Method, Kernels	Optimal cost function value	$\delta_t$	$\sigma_t$	$\sigma^2$	$\alpha_{se}$	$c_t$	$h_t$	$\theta_t$
MLM, (40)+(44)	1.8876e+06	4.6377e+03	5000	173.0302	0.1303	0.3834	109.4592	0.0663
MLM, (40)+(46)	2.0311e+06	361.7502	2.2946	501.5663	8.0899	—	—	—
GCV, (40)+(44)	242.3078	9.3410e+03	5000	173.0302	0.1373	0.5345	93.4106	0.1136
GCV, (40)+(46)	508.7996	9.0237e+08	1.1407e+07	501.5655	2.5726	—	—	—
SURE, (40)+(44)	9.8970e+07	9.2773e+03	5000	173.0302	0.1376	0.5346	92.7727	0.0562
SURE, (40)+(46)	2.2675e+08	5.0370e+08	6.8254e+06	501.5655	2.8876	—	—	—

Table 4

The average prediction fits (48) for the GHCN temperature data using two kernel combinations, where the hyper-parameters are estimated by MLM, GCV and SURE methods, respectively. The values in parentheses are the corresponding smallest prediction fits.

Kernel	(40)+(44)	(40)+(46)
MLM	60.34 (44.12)	-62.85(-146.20)
GCV	60.30 (35.97)	-57.48(-232.58)
SURE	52.40 (22.20)	-53.57(-226.53)

First, we recall from e.g., [18], that the subsystem at the  $i$ th location  $p_i$  with  $i = 1, \dots, M$  of a spatially-distributed system can be described by the following ARX model

$$\mathcal{A}_i(q_p, q_t) f(p_i, t_j) = \mathcal{B}_i(q_p, q_t) u(p_i, t_j), \quad (50)$$

$$t_j = jT_s, j = 1, 2, \dots, N,$$

where  $q_p$  and  $q_t$  are the forward spatial and temporal shift operators, respectively, i.e.,  $q_p q_t^{-1} f(p_i, t_j) = f(p_{i+1}, t_{j-1})$ ,  $f(p_i, t_j) \in \mathbb{R}$  and  $u(p_i, t_j) \in \mathbb{R}$  are the output and input at the  $i$ th location  $p_i$  and  $j$ th time instant  $t_j$ , respectively, and

$$\mathcal{A}_i(q_p, q_t) = 1 + \sum_{k_i=1}^{n_a} \sum_{k_j=0}^M a_i(p_{k_j}, t_{k_i}) q_p^{-k_j} q_t^{-k_i}, \quad (51)$$

$$\mathcal{B}_i(q_p, q_t) = \sum_{l_i=1}^{n_b} \sum_{l_j=0}^M b_i(p_{l_j}, t_{l_i}) q_p^{-l_j} q_t^{-l_i}, \quad (52)$$

with  $a_i(p_{k_j}, t_{k_i}), b_i(p_{l_j}, t_{l_i}) \in \mathbb{R}$ ,  $n_a, n_b \in \mathbb{N}$  and  $p_0$  a null position. Then, we consider a special case of (50) with

$$\mathcal{A}_i(q_p, q_t) = 1, \quad (53)$$

$$\mathcal{B}_i(q_p, q_t) = \sum_{k=1}^{n_b} b_i(p_i, t_k) q_t^{-k}, u(p_i, t_j) = u(t_j),$$

which is equivalent to assume that the subsystem at the  $i$ th location  $p_i$  with  $i = 1, \dots, M$ , has a finite impulse response (FIR) model. The FIR parameters  $b_i(p_i, t_k)$ ,  $k = 1, \dots, n_b$  only depend on  $p_i$  and moreover, assumed

to be smooth functions of  $p_i$ . In this case, the output  $f(p_i, t_j)$  takes the form of

$$f(p_i, t_j) = \sum_{k=1}^{n_b} b_i(p_i, t_k) u(t_{j-k}), \quad (54)$$

which plays the role as the spatial-temporal function  $f(p_i, t_j)$  in (1).

In what follows, we study the identification of spatially-distributed system (50) with (53), i.e., the estimation of the FIRs  $\{b(p_i, t_j)\}_{i=1, j=1}^{M, n_b}$  of  $M$  spatially-distributed subsystems as well as possible based on the training data  $\{y_{i,j}, u(t_j)\}_{i=1, j=1}^{M, N}$  by using the Gaussian process regression approach in this paper. For comparison, we also consider the estimation of the FIR  $\{b(p_i, t_j)\}_{j=1}^{n_b}$  of the  $i$ th subsystem based on  $\{y_{i,j}, u(t_j)\}_{j=1}^N$  separately by neglecting the spatial interconnections between  $M$  subsystems and by using the approach in [12]. These two approaches are denoted by the ‘‘spatial-temporal’’ and ‘‘temporal’’ approaches in the following, respectively.

#### 4.3.1 Test Spatially-distributed Systems

We first generate a 30th order discrete time system using the procedure in [12] with 5 poles with the largest modulus lying in  $[0.8, 0.9]$  and one pole with the 6th largest modulus smaller than 0.75. For convenience, we let  $\mathbf{p}_k$ , whose real part is  $\mathbf{a}_k$  and the imaginary part  $\mathbf{b}_k$ , denote the pole of this system with the  $k$ th largest modulus. Then we generate  $M = 500$  new test systems by keeping zeros and poles of this system unchanged except the 5 poles with the largest modulus, i.e.,  $\{\mathbf{p}_k\}_{k=1}^5$  and then. Note that if there exists an unpaired non-real pole in  $\{\mathbf{p}_k\}_{k=1}^5$ , e.g.,  $\mathbf{a}_k + \mathbf{b}_k i$  with  $\mathbf{b}_k \neq 0$  is included but  $\mathbf{a}_k - \mathbf{b}_k i$  is not, we will regenerate the original system until  $\{\mathbf{p}_k\}_{k=1}^5$  include either real poles or complex conjugate pairs of poles. For the  $i$ th new test system with  $i = 1, \dots, 500$ ,  $\{\mathbf{p}_k\}_{k=1}^5$  are modified as  $\{\mathbf{p}_{k,i}\}_{k=1}^5$  with the real part  $\mathbf{a}_{k,i}$  and the imaginary part  $\mathbf{b}_{k,i}$  as follows,

- for real pole  $\mathbf{p}_k$ , the modified  $\mathbf{p}_{k,i} = \mathbf{a}_{k,i}$  is uniformly distributed in  $[\mathbf{a}_k - 0.05, \mathbf{a}_k + 0.05]$  and  $\mathbf{b}_{k,i} = 0$ ;

Table 5

The hyper-parameter and the corresponding optimal values of the cost functions for the GHCN temperature data using two kernel combinations: (40)+(44) and (40)+(46), and three hyper-parameter estimation methods, respectively.

Method, Kernels	Optimal cost function value	$\delta_t$	$\sigma_t$	$\sigma^2$	$\alpha_{se}$	$c_t$	$h_t$	$\theta_t$
MLM, (40)+(44)	5.4912e+07	585.9242	5000	2.6159	984.0928	0.2867	5.8592	1.9000
MLM, (40)+(46)	5.7639e+07	626.4846	235.5695	1.1125	91.8811	-	-	-
GCV, (40)+(44)	2.1015	1.8304e+04	5000	2.6159	19.6885	0.6481	1.7664e+03	8.3648
GCV, (40)+(46)	2.4189	2.4023e+04	6.1204e+03	1.1125	47.8679	-	-	-
SURE, (40)+(44)	7.6828e+07	5.7973e+03	5000	2.6159	345.0477	0.4735	579.7336	50.3385
SURE, (40)+(46)	4.4955e+07	1.1622e+04	4.8269e+03	1.1125	64.0244	-	-	-

- for complex conjugate pair of poles  $\mathbf{p}_k = \mathbf{a}_k \pm \mathbf{b}_k i$ , the modified complex conjugate pair of poles are  $\mathbf{a}_{k,i} \pm \mathbf{b}_{k,i} i$ , where  $(\mathbf{a}_{k,i}, \mathbf{b}_{k,i})$  is uniformly distributed in the circle with the center  $(\mathbf{a}_k, \mathbf{b}_k)$  and radius 0.05.

Now we obtain 500 spatially-distributed test subsystems and for the  $i$ th subsystem, the corresponding location is  $p_i = [\mathbf{a}_{1,i}, \mathbf{b}_{1,i}, \dots, \mathbf{a}_{5,i}, \mathbf{b}_{5,i}]^T \in \mathbb{R}^{10}$ .

#### 4.3.2 Test Data Sets

We choose the test input signal  $u(t_j) = e^{-\alpha t_j} \sin(\omega_0 t_j)$  with  $\alpha = 10^{-2}$  and  $\omega_0 = \pi/8$ , whose state-space model is in the form of

$$\begin{aligned} \tilde{z}_{i,j+1} &= \tilde{E} \tilde{z}_{i,j} + \tilde{F} \delta_j, \quad \tilde{z}_{i,0} = 0 \in \mathbb{R}^2, \\ u(t_j) &= \tilde{H} \tilde{z}_{i,j}, \quad j = 0, 1, \dots, N, \end{aligned} \quad (55)$$

where  $\tilde{z}_{i,j} \in \mathbb{R}^2$ ,  $\delta_j$  denotes the impulsive input, i.e.,  $\delta_j = 1$  for  $j = 0$  and  $\delta_j = 0$  for  $j = 1, \dots$ , and

$$\tilde{E} = \begin{bmatrix} 2e^{-\alpha} \cos(\omega_0) & -e^{-2\alpha} \\ 1 & 0 \end{bmatrix}, \quad (56a)$$

$$\tilde{F} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \quad \tilde{H} = \begin{bmatrix} e^{-\alpha} \sin(\omega_0) \\ 0 \end{bmatrix}^T. \quad (56b)$$

Then for  $i = 1, \dots, M$ , we simulate the  $i$ th test subsystem with the test input signal to get the noise-free output  $f(p_i, t_j)$  and then corrupt it with an additive measurement noise  $v_{i,j}$ , which follows a Gaussian distribution with zero mean and variance  $\sigma^2$ , leading to a data record with 400 pairs of input and measurement output data  $\{y_{i,j}, u(t_j)\}_{j=1}^{400}$ . The average signal-to-noise ratio (SNR) of 500 test subsystems is 1, where the SNR of each test subsystem is defined as the ratio between the variance of the noise-free output  $f(p_i, t_j)$  and that of the measurement noise  $v_{i,j}$ . In this way, the generated data sets contain 500 data records, each with 400 pairs of input and measurement output data, i.e.,  $\{y_{i,j}, u(t_j)\}_{i=1, j=1}^{500, 400}$ .

#### 4.3.3 Choice of Kernels

For the ‘‘spatial-temporal’’ approach, the spatial kernel  $k_s(p_i, p_{i'}; \alpha_s)$  and the temporal kernel  $k_t(t_j, t_{j'}; \alpha_t)$  in (4)

are chosen to be the SE kernel (40), and the following one, respectively,

$$k_t(t_j, t_{j'}; \alpha_t) = \kappa(t_j, t_{j'}; \alpha_t) \sum_{k=1}^{n_b} \sum_{k'=1}^{n_b} u(t_{j-k}) u(t_{j'-k'}), \quad (57)$$

where  $\kappa(t_j, t_{j'}; \alpha_t) : \mathbb{R}_+ \times \mathbb{R}_+ \rightarrow \mathbb{R}$  is the diagonal correlated (DC) kernel in [12], i.e.,

$$\begin{aligned} \kappa(t_j, t_{j'}; \alpha_t) &= \delta_t \lambda_t^{(t_j+t_{j'})/2} \frac{|t_j-t_{j'}|}{\rho_t}, \\ \alpha_t &= [\delta_t, \lambda_t, \rho_t] \in \Omega = \{\delta_t \geq 0, \lambda_t \in [0, 1], |\rho_t| \leq 1\}. \end{aligned} \quad (58)$$

Noting the state-space model realization of the DC kernel (58) in [9] and (55), it can be shown that the state-space model realization of the spatial-temporal kernel (4) with (40) as the spatial kernel and (57) as the temporal kernel takes the form of (19) by replacing (19a) with

$$s_{j+1} = F s_j + G_j w_j, \quad s_1 \sim \mathcal{N}(0, I_M \otimes \Sigma_1),$$

and using

$$r = 3, \quad F = \begin{bmatrix} I_M \otimes (\lambda_t^{1/2} \rho_t) & 0 \in \mathbb{R}^{M \times 2M} \\ I_M \otimes ((1 - \rho_t^2)^{1/2} \tilde{F}) & I_M \otimes \tilde{E} \end{bmatrix},$$

$$G_j = \begin{bmatrix} (I_M \otimes \lambda_t^{1/2}) \delta_t^{1/2} \lambda_t^{t_j/2} \\ 0 \in \mathbb{R}^{2M \times M} \end{bmatrix}, \quad H = \begin{bmatrix} 0 \in \mathbb{R}^{M \times M} \\ (I_M \otimes \tilde{H}) \end{bmatrix}^T,$$

$$\Sigma_1 = \begin{bmatrix} I_M \otimes (\delta_t / (1 - \rho_t^2)) & 0 \in \mathbb{R}^{M \times 2M} \\ 0 \in \mathbb{R}^{2M \times M} & 0 \in \mathbb{R}^{2M \times 2M} \end{bmatrix},$$

where  $\tilde{E}, \tilde{F}$  and  $\tilde{H}$  are given in (56). Note that for the ‘‘temporal’’ approach, we only apply the DC kernel (58).

#### 4.3.4 Hyper-parameter estimation and Impulse Response Estimation

For the ‘‘spatial-temporal’’ approach, we use the MLM method (5) to estimate  $\alpha = [\alpha_s^T, \alpha_t^T, \sigma^2]^T$  and apply the same strategy as stated in Section 4.2.3 for finding a ‘‘good’’ local minimum. With the estimated hyper-parameter, we further run the Kalman filter and

smoother to obtain the estimates of  $\{b_i(p_i, t_j)\}_{i=1, j=1}^{M, n_b}$ , denoted as  $\{\hat{b}_{i,j|N}\}_{i=1, j=1}^{M, n_b}$ , where for  $j = 1, \dots, n_b$ ,

$$\hat{b}_{j|N} = \left[ \hat{b}_{1,j|N} \cdots \hat{b}_{M,j|N} \right]^T \in \mathbb{R}^M, \quad (59)$$

can be obtained by

$$\hat{b}_{j|N} = \Lambda D^{1/2} (I_M \otimes (1 - \rho_t^2)^{1/2}) [\hat{x}_{j|N}]_{1:M}. \quad (60)$$

Here  $\Lambda$  and  $D$  are defined in (20), and  $[\hat{x}_{j|N}]_{1:M}$  denotes a vector containing the first  $M$  elements of  $\hat{x}_{j|N}$  in (30c).

For the ‘‘temporal’’ approach, we use the MLM method (5) to estimate  $\alpha = [\alpha_t^T, \sigma^2]^T$  and then with the estimated hyper-parameter, we calculate  $\hat{b}_i = [\hat{b}_{i,1|N}, \dots, \hat{b}_{i,n_b|N}]^T \in \mathbb{R}^{n_b}$  for the  $i$ th system with  $i = 1, \dots, M$ , where the implementation [11] is used.

To evaluate the estimation performance of  $\{\hat{b}_{i,j|N}\}_{i=1, j=1}^{M, n_b}$ , for the  $i$ th system with  $i = 1, \dots, M$ , we let

$$b_i^0 = \left[ b_{i,1}^0 \cdots b_{i,n_b}^0 \right]^T,$$

denote the true value of  $[b_i(p_i, t_1), \dots, b_i(p_i, t_{n_b})]^T$ , and then define the measure of fit, e.g., [20],

$$\text{fit}_i^b = 100 \times \left( 1 - \frac{\|\hat{b}_i - b_i^0\|_2}{\|\hat{b}_i - \bar{b}_i^0\|_2} \right), \quad \bar{b}_i^0 = \frac{1}{n_b} \sum_{j=1}^{n_b} b_{i,j}^0.$$

The average estimation fit of  $\{\hat{b}_i\}_{i=1}^M$  is defined as

$$\overline{\text{fit}}^b = \frac{1}{M} \sum_{i=1}^M \text{fit}_i^b. \quad (61)$$

#### 4.3.5 Simulation Results and Findings

In the simulation, we choose the FIR order  $n_b = 125$ . The average estimation fits of  $\{\hat{b}_i\}_{i=1}^M$  of the ‘‘spatial-temporal’’ and ‘‘temporal’’ approaches in Table 6 show that the ‘‘spatial-temporal’’ approach gives much better estimation performance than the ‘‘temporal’’ approach. This observation indicates that exploring the spatial interconnections among subsystems is beneficial for the identification of spatially-distributed system.

Table 6  
Average estimation fits of  $M$  spatially-distributed systems

Approach	‘‘spatial-temporal’’	‘‘temporal’’
$\overline{\text{fit}}^b$ (61)	77.27	8.78

## 5 Conclusion

In this paper, we proposed an efficient implementation with computational complexity  $\mathcal{O}(M^3 + NM^2)$ , for spatial-temporal Gaussian process regression by exploring the Kronecker structure of its state-space model realization, where  $N$  and  $M$  are the numbers of time instants and locations, respectively. The proposed implementation has been illustrated over applications in weather data prediction and spatially-distributed system identification. For the weather prediction, the design kernel is shown to give better prediction performance than the one in [27] and for the spatially-distributed system identification, the benefit of exploring the spatial interconnections among subsystems is confirmed.

## Appendix A

This appendix contains the proofs of all theoretical results and the derivations of state-space model of (42).

### A.1 Proof of Proposition 1

According to (32a)-(32b) in Lemma 1 and (33a), we have

$$\Theta = \begin{bmatrix} l_1 \\ l_2 - b_{2,1} l_1 \\ \vdots \\ l_N - \sum_{i=1}^{N-1} b_{N,i} l_i \end{bmatrix} = \Gamma L, \quad (A.1)$$

where  $\Gamma$  is defined in (35). Then, inserting (34a) into (33b), it follows that  $\Psi = \text{COV}[\Gamma L, \Gamma L] = \Gamma \text{COV}[L, L] \Gamma^T$ , which leads to (34b) using (25b). Combining (28b), (33b) and (32c), we can obtain (34c).

### A.2 Proof of Proposition 2

First, note that the computation of the cost function of (5) depends on that of  $\log |\Sigma(\alpha)|$  and  $Y^T \Sigma^{-1}(\alpha) Y$ . Then following the idea of [5], where the computation of the generalized cross validation filter is discussed, and using (25) and Proposition 1,  $\log |\Sigma(\alpha)|$  and  $Y^T \Sigma^{-1}(\alpha) Y$  can be computed as follows

$$\begin{aligned} & \log |\Sigma(\alpha)| \\ &= \log |(I_N \otimes \Lambda) \bar{\Sigma}(\alpha) (I_N \otimes \Lambda^T)| \quad (A.2a) \\ &= \log |(I_N \otimes \Lambda) (I_N \otimes \Lambda^T)| + \log |\bar{\Sigma}(\alpha)| = \log |\Psi| \\ & \quad Y^T \Sigma^{-1}(\alpha) Y \\ &= L^T (I_N \otimes \Lambda^T) (I_N \otimes \Lambda) \bar{\Sigma}(\alpha)^{-1} (I_N \otimes \Lambda^T) (I_N \otimes \Lambda) L \\ &= (\Gamma^{-1} \Theta)^T (\Gamma^T \Psi \Gamma)^{-1} (\Gamma^{-1} \Theta) = \Theta^T \Psi^{-1} \Theta, \quad (A.2b) \end{aligned}$$

where the first steps of both (A.2a) and (A.2b) are derived from (25), and the second step of (A.2b) is derived

from (34a) and (34b). Then using (33a) and (34c), we can obtain (36).

### A.3 Proof of Proposition 3

As shown in (6) and (7), the computation of the cost functions of the GCV and SURE methods depends on that of  $\delta$  and  $S$  defined in (8f) and (8g), respectively. For convenience, we let  $\gamma = \sigma^2$ . We first rewrite (8a) as

$$\gamma \Sigma(\alpha)^{-1} = I_{NM} - [K_t(\alpha_t) \otimes K_s(\alpha_s)] \Sigma(\alpha)^{-1}. \quad (\text{A.3})$$

Following the discussions in [5], we can represent  $\delta$  and  $S$  as functions of  $\log |\Sigma(\alpha)|$  and  $Y^T \Sigma^{-1}(\alpha) Y$ , respectively,

$$\begin{aligned} & \gamma \frac{\partial \log |\Sigma(\alpha)|}{\partial \gamma} \\ &= \gamma \text{trace}(\Sigma(\alpha)^{-1} \frac{\partial \Sigma(\alpha)}{\partial \gamma}) = \gamma \text{trace}(\Sigma(\alpha)^{-1}) \\ &= \text{trace} \{ I_{NM} - [K_t(\alpha_t) \otimes K_s(\alpha_s)] \Sigma(\alpha)^{-1} \} \\ &= NM - \text{trace} \{ [K_t(\alpha_t) \otimes K_s(\alpha_s)] \Sigma(\alpha)^{-1} \} \\ &= NM - \delta, \quad (\text{A.4}) \\ &- \gamma^2 \frac{\partial Y^T \Sigma(\alpha)^{-1} Y}{\partial \gamma} \\ &= \gamma^2 Y^T \Sigma^{-1} \frac{\partial \Sigma(\alpha)}{\partial \gamma} \Sigma^{-1} Y = \gamma^2 Y^T \Sigma^{-2} Y \\ &= \gamma^2 Y^T \{ I_{NM} - [K_t(\alpha_t) \otimes K_s(\alpha_s)] \Sigma(\alpha)^{-1} \}^T \\ &\quad \{ I_{NM} - [K_t(\alpha_t) \otimes K_s(\alpha_s)] \Sigma(\alpha)^{-1} \} Y \\ &= \|\hat{Y} - Y\|_2^2 = S. \quad (\text{A.5}) \end{aligned}$$

Then by using (36),  $S$  and  $\delta$  can be computed as follows

$$\begin{aligned} S &= -\gamma^2 \frac{\partial Y^T \Sigma(\alpha)^{-1} Y}{\partial \gamma} = -\gamma^2 \sum_{j=1}^N \frac{\partial \bar{e}_j^T \bar{E}_j^{-1} \bar{e}_j}{\partial \gamma}, \\ \delta &= MN - \gamma \frac{\partial \log |\Sigma(\alpha)|}{\partial \gamma} = MN - \gamma \sum_{j=1}^N \frac{\partial \log |\bar{E}_j|}{\partial \gamma}. \quad (\text{A.6}) \end{aligned}$$

Now we define

$$\bar{\zeta}_{j|j-1} = \frac{\partial \hat{x}_{j|j-1}}{\partial \gamma}, \bar{P}_{j|j-1} = \frac{\partial \bar{\Sigma}_{j|j-1}}{\partial \gamma}, \quad (\text{A.7})$$

and then  $\partial \bar{e}_j^T \bar{E}_j^{-1} \bar{e}_j / \partial \gamma$  and  $\partial \log |\bar{E}_j| / \partial \gamma$  in (A.6) can be further expressed as

$$\begin{aligned} - \frac{\partial \bar{e}_j^T \bar{E}_j^{-1} \bar{e}_j}{\partial \gamma} &= \bar{e}_j^T \bar{E}_j^{-1} (\bar{H} \bar{P}_{j|j-1} \bar{H}^T + I_M) \bar{E}_j^{-1} \bar{e}_j, \\ &\quad + 2 \bar{\zeta}_{j|j-1}^T \bar{H} \bar{E}_j^{-1} \bar{e}_j, \quad (\text{A.8a}) \end{aligned}$$

$$\frac{\partial \log |\bar{E}_j|}{\partial \gamma} = \text{trace} [\bar{E}_j^{-1} (\bar{H} \bar{P}_{j|j-1} \bar{H}^T + I_M)]. \quad (\text{A.8b})$$

Combining (A.6) with (A.8), we can obtain (37). Moreover, inserting (28d) and (28f) into (A.7), and (28e) and (28g) into (A.7), we can compute  $\bar{\zeta}_{j|j-1}$  and  $\bar{P}_{j|j-1}$  recursively as shown in Proposition 3.

### A.4 Proof of Theorem 1

As shown in Algorithm 1, the proposed implementation consists of three steps, and in what follows, we will study their computational complexities, respectively:

- 1) *Computational complexity of Step 1:* We first calculate the SVD of  $K_s \in \mathbb{R}^{M \times M}$  in (20) and its computational complexity is  $\mathcal{O}(M^3)$ . Then the computational complexities of (24a) and (25a) are  $\mathcal{O}(Mr)$  and  $\mathcal{O}(NM^2)$ , respectively. Hence, this step has the computational complexity  $\mathcal{O}(M^3 + NM^2)$ .
- 2) *Computational complexity of Step 2:* Since the evaluation of the cost functions of three hyper-parameter estimation methods all rely on Kalman filter, we first consider the computational complexity of the Kalman filter (28) and then that of three hyper-parameter estimation methods, respectively.

- (a) *Computational complexity of Kalman filter (28):* To show the computational complexity of the Kalman filter, we first use induction to show that, for  $j = 1, \dots, N$ ,  $\bar{E}_j \in \mathbb{R}^{M \times M}$  and  $\bar{\Sigma}_{j|j-1} \in \mathbb{R}^{Mr \times Mr}$  are diagonal and block diagonal matrices, respectively. It consists of two steps.

Our first step is to prove that  $\bar{E}_1$  and  $\bar{\Sigma}_{1|0}$  are diagonal and block diagonal matrices, respectively. For  $\bar{\Sigma}_{1|0}$ , inserting (24a), (29) and  $\bar{\Sigma}_{0|0} = \mathbb{E}[(x_0 - \mathbb{E}(x_0))(x_0 - \mathbb{E}(x_0))^T] = I_M \otimes \Sigma_0$ , where we apply (23a) and (26), into (28g), we have

$$\begin{aligned} \bar{\Sigma}_{1|0} &= I_M \otimes (F_D \Sigma_0 F_D^T + G_D G_D^T) \\ &= \text{blkdiag}(\bar{\Sigma}_{1,1}, \dots, \bar{\Sigma}_{1,M}), \quad (\text{A.9}) \end{aligned}$$

where  $\bar{\Sigma}_{1,i} = F_D \Sigma_0 F_D^T + G_D G_D^T$  for  $i = 1, \dots, M$ . For  $\bar{E}_1$ , we insert (24a) and (A.9) into (28b) to obtain

$$\begin{aligned} \bar{E}_1 &= H_D \bar{\Sigma}_{1,1} H_D^T D + \sigma^2 I_M \\ &= \text{diag}(\bar{E}_{1,1}, \dots, \bar{E}_{1,M}), \quad (\text{A.10}) \end{aligned}$$

where  $\bar{E}_{1,i} = [D]_{ii} H_D \bar{\Sigma}_{1,1} H_D^T + \sigma^2$  for  $i = 1, \dots, M$ .

Our second step is to show that for  $j = 1, \dots, N-1$ , if we assume  $\bar{E}_j$  and  $\bar{\Sigma}_{j|j-1}$  are diagonal and block diagonal matrices, respectively, then we can show that  $\bar{E}_{j+1}$  and  $\bar{\Sigma}_{j+1|j}$  are diagonal and block diagonal matrices, respectively. For convenience, we define that

$$\bar{E}_j = \text{diag}(\bar{E}_{j,1}, \dots, \bar{E}_{j,M}), \quad (\text{A.11a})$$

$$\bar{\Sigma}_{j|j-1} = \text{blkdiag}(\bar{\Sigma}_{j,1}, \dots, \bar{\Sigma}_{j,M}), \quad (\text{A.11b})$$

where  $\bar{E}_{j,i} \in \mathbb{R}$  and  $\bar{\Sigma}_{j,i} \in \mathbb{R}^{r \times r}$  for  $i = 1, \dots, M$ . Combining (28g) and (28e), we have

$$\begin{aligned} \bar{\Sigma}_{j+1|j} &= \bar{F} \bar{\Sigma}_{j|j-1} \bar{F}^T + Q \\ &\quad - \bar{F} \bar{\Sigma}_{j|j-1} \bar{H}^T \bar{E}_j^{-1} \bar{H} \bar{\Sigma}_{j|j-1} \bar{F}^T \\ &= \text{blkdiag}(\bar{\Sigma}_{j+1,1}, \dots, \bar{\Sigma}_{j+1,M}), \end{aligned} \quad (\text{A.12})$$

$$\begin{aligned} \bar{\Sigma}_{j+1,i} &= F_D \bar{\Sigma}_{j,i} F_D^T + G_D G_D^T \\ &\quad + ([D]_{ii} / \bar{E}_{j,i}) F_D \bar{\Sigma}_{j,i} H_D^T H_D \bar{\Sigma}_{j,i} F_D^T, \end{aligned} \quad (\text{A.13})$$

where  $i = 1, \dots, M$ , and we apply (24a), (29), (A.11) and the fact that  $\bar{H}^T \bar{E}_j^{-1} \bar{H} = (I_M \otimes H_D^T) [(D^{1/2} \bar{E}_j^{-1} D^{1/2}) \otimes 1] (I_M \otimes H_D) = (D^{1/2} \bar{E}_j^{-1} D^{1/2}) \otimes (H_D^T H_D)$  is a block diagonal matrix with  $i$ th block being  $([D]_{ii} / \bar{E}_{j,i}) H_D^T H_D \in \mathbb{R}^{r \times r}$ . Then for  $\bar{E}_{j+1}$ , we use (24a) and (A.12) to obtain

$$\begin{aligned} \bar{E}_{j+1} &= \bar{H} \bar{\Sigma}_{j+1|j} \bar{H}^T + \sigma^2 I_M \\ &= \text{diag}(\bar{E}_{j+1,1}, \dots, \bar{E}_{j+1,M}), \end{aligned} \quad (\text{A.14})$$

where  $\bar{E}_{j+1,i} = [D]_{ii} H_D \bar{\Sigma}_{j+1,i} H_D^T + \sigma^2$  for  $i = 1, \dots, M$ .

Hence for  $j = 1, \dots, N$ , it is clear that  $\bar{E}_j \in \mathbb{R}^{M \times M}$  and  $\bar{\Sigma}_{j|j-1} \in \mathbb{R}^{Mr \times Mr}$  are diagonal and block diagonal matrices, respectively. It follows that  $\bar{\Sigma}_{j|j} \in \mathbb{R}^{Mr \times Mr}$  in (28e) is also a block diagonal matrix due to that  $\bar{H}^T \bar{E}_j^{-1} \bar{H}$  is a block diagonal matrix.

Then, according to the properties of the Kronecker product and the matrix multiplication, for each  $j = 1, \dots, N$ , the computational complexities of (28b), (28c) and (28d), and (28e) and (28g) are  $\mathcal{O}(Mr^2)$  and  $\mathcal{O}(Mr^3)$ , respectively. Thus the computational complexity of (28) is  $\mathcal{O}(NM)$ .

- (b) *Computational complexity of the computation of the cost function of the MLM method:* As shown in Proposition 2, to calculate the cost function of the MLM method (5), we first calculate  $\bar{e}_j$  and  $\bar{E}_j$  as shown in (28) for  $j = 1, \dots, N$  and the computational complexity is  $\mathcal{O}(NM r^3)$ . Then we calculate (36) and (5), whose computational complexity is  $\mathcal{O}(NM)$ . Therefore, the computational complexity of the cost function of the MLM method is  $\mathcal{O}(NM)$ .
- (c) *Computational complexities of the computation of the cost functions of the GCV and SURE methods:* For the GCV method (6) and the SURE method (7), as shown in Proposition 3, since  $\bar{E}_j$  is a diagonal matrix, and  $\bar{\Sigma}_{j|j-1}$ ,  $\bar{F}$  and  $\bar{H}$  are block diagonal matrices, the computational complexity of  $\bar{\Sigma}_{j|j-1} \in \mathbb{R}^{Mr}$  and  $\bar{P}_{j|j-1} \in \mathbb{R}^{Mr \times Mr}$  with  $j = 1, \dots, N$  in (38) and (39) are  $\mathcal{O}(Mr^3)$ . Therefore, the computational complexities of the

cost functions of the GCV and the SURE methods are both  $\mathcal{O}(NM)$ .

- 3) *Computational complexity of Step 3:* We discuss the computational complexities of the Kalman smoother (30) and predictor (31), respectively.
- (a) *Computational complexity of Kalman smoother:* Since  $\bar{\Sigma}_{j|j} \in \mathbb{R}^{Mr \times Mr}$ ,  $j = N-1, \dots, 1$  are block diagonal matrices,  $\bar{J}_j$ ,  $j = N-1, \dots, 1$  are also block diagonal matrices. Hence, for each  $j$ , the computational complexities of (30a), and (30b), (30c) and (30d) are  $\mathcal{O}(Mr^2)$  and  $\mathcal{O}(Mr^3)$ , respectively. Finally, since the output transform in (30e) has computational complexity  $\mathcal{O}(M^2 r)$ , the computational complexity of (30) is  $\mathcal{O}(NM^2)$ .
- (b) *Computational complexity of Kalman predictor:* For each  $j$ , the computational complexities of (31b) and (31c) are  $\mathcal{O}(Mr^3)$  and  $\mathcal{O}(M^2 r)$ , respectively. Thus the computational complexity of (31) is  $\mathcal{O}(NM^2)$ .

Hence, the proof of Theorem 1 is complete.

#### A.5 State-space Model Realization of (42)

The kernel (42) can be divided into three parts:

$$\begin{aligned} & k_{\text{Te2}}(\tau; \delta_t, c_t) k_{\text{EXP}}(\tau) \quad (\text{A.15}) \\ &= \underbrace{\delta_t [1 - c_t + \frac{3}{4} c_t^2] \exp\left(-\frac{|\tau|}{\sigma_t}\right)}_{(\text{A.15a})} \\ &\quad + \underbrace{(c_t - c_t^2) \cos(2\pi \mathbf{f} |\tau|) \exp\left(-\frac{|\tau|}{\sigma_t}\right)}_{(\text{A.15b})} \\ &\quad + \underbrace{\frac{c_t^2}{4} \cos(4\pi \mathbf{f} |\tau|) \exp\left(-\frac{|\tau|}{\sigma_t}\right)}_{(\text{A.15c})}, \end{aligned}$$

where (A.15a) is an exponential kernel, and (A.15b) and (A.15c) are periodic kernels with different periods. To obtain the state-space model of (42), we derive below the state-space models of these kernels, respectively.

Firstly, we denote the PSD of (A.15a) as  $\Phi_{\text{EXP}}(\omega)$ , which can be obtained using (10) as follows

$$\Phi_{\text{EXP}}(\omega) = \frac{\delta_t (1 - c_t + \frac{3}{4} c_t^2) (1 - e^{-2\beta_t})}{(e^{-i\omega} - e^{-\beta_t}) (e^{i\omega} - e^{-\beta_t})}, \quad (\text{A.16})$$

where  $\beta_t = \frac{1}{\sigma_t}$ . According to Assumption 1 and (11), we can obtain the transfer function of (A.16) in the form

$$W_{\text{EXP}}(e^{i\omega}) = \frac{\sqrt{\delta_t (1 - c_t + \frac{3}{4} c_t^2) (1 - e^{-2\beta_t})}}{e^{i\omega} - e^{-\beta_t}}. \quad (\text{A.17})$$

Then we can derive the corresponding state-space model using the realization theory in [8] as follows

$$s_{1,j+1} = F_1 s_{1,j} + G_1 w_{1,j}, j = 1, \dots \quad (\text{A.18a})$$

$$y_{1,j} = H_{1,j} s_{1,j}, \quad (\text{A.18b})$$

where  $s_{1,j} \in \mathbb{R}$ ,  $F_1 = e^{-\beta t}$ ,  $G_1 = 1$  and  $H_{1,j} = \sqrt{\delta_t(1 - c_t + \frac{3}{4}c_t^2)(1 - e^{-2\beta t})}$  and  $w_{1,j} \in \mathbb{R}$  is white Gaussian noise with zero mean and unit variance.

Secondly, we denote the PSD of (A.15b) as  $\Phi_{\text{PD1}}(\omega)$  and it can be derived using (10) as follows

$$\Phi_{\text{PD1}}(\omega) = \delta_t(c_t - c_t^2) \quad (\text{A.19})$$

$$\left[ \frac{1}{2} (e^{3\beta t} - e^{\beta t}) (e^{-i\varrho_1} + e^{i\varrho_1}) (e^{-i\omega} + e^{i\omega}) + 1 - e^{4\beta t} \right] / \left[ (e^{\beta t} - e^{-i\varrho_1 - i\omega}) (e^{\beta t} - e^{i\varrho_1 - i\omega}) (e^{\beta t} - e^{i\omega - i\varrho_1}) (e^{\beta t} - e^{i\varrho_1 + i\omega}) \right],$$

where  $\varrho_1 = 2\pi f$ . According to Assumption 1 and (11), we consider the transfer function of (A.19) in the form

$$W(e^{i\omega}) = \sqrt{\delta_t(c_t - c_t^2)} \frac{n_1 e^{i\omega} + n_2}{e^{2i\omega} - d_2 e^{i\omega} - d_1}, \quad (\text{A.20})$$

with  $d_1 = -e^{-2\beta t}$ ,  $d_2 = e^{-\beta t}(e^{i\varrho_1} + e^{-i\varrho_1})$ ,

$$A = n_1^2 + n_2^2 = 1 - e^{4\beta t}, \quad (\text{A.21a})$$

$$B = n_1 n_2 = \frac{1}{2} (e^{3\beta t} - e^{\beta t}) (e^{-i\varrho_1} + e^{i\varrho_1}). \quad (\text{A.21b})$$

Therefore, we can derive  $n_1$  and  $n_2$  by solving

$$n_1^4 - A n_1^2 + B^2 = 0. \quad (\text{A.22})$$

With (A.20),  $n_1$  and  $n_2$ , we can derive the state-space model in controllable canonical form of (A.15b) using the realization theory in [8] as follows

$$s_{2,j+1} = F_2 s_{2,j} + G_2 w_{2,j}, j = 1, \dots \quad (\text{A.23a})$$

$$y_{2,j} = H_{2,j} s_{2,j}, \quad (\text{A.23b})$$

where  $s_{2,j} \in \mathbb{R}^2$ ,  $w_{2,j} \in \mathbb{R}^2$  is white Gaussian noise with zero mean and covariance matrix  $I_2$ , and

$$F_2 = \begin{bmatrix} 0 & 1 \\ d_1 & d_2 \end{bmatrix}, G_2 = \begin{bmatrix} 0 \\ 1 \end{bmatrix}, H_{2,j} = \sqrt{\delta_t(c_t - c_t^2)} \begin{bmatrix} n_2 & n_1 \end{bmatrix}.$$

Moreover, it is necessary to check if the  $n_1$  and  $n_2$  guarantees that the zeros of (A.20) are inside the unit circle and the largest eigenvalue of the corresponding system matrices  $F_2$  should be less than 1.

The PSD of (A.15c) can be obtained by in a similar way as (A.15b) by replacing the amplitude with  $\delta_t \frac{c_t^2}{4}$  and letting  $\varrho_1 = 4\pi f$  in (A.19), respectively. Let the state-space model of (A.15c) in controllable canonical form be represented as

$$s_{3,j+1} = F_3 s_{3,j} + G_3 w_{3,j}, j = 1, \dots \quad (\text{A.24a})$$

$$y_{3,j} = H_{3,j} s_{3,j}, \quad (\text{A.24b})$$

where  $s_{3,j} \in \mathbb{R}^2$ ,  $F_3 \in \mathbb{R}^{2 \times 2}$ ,  $G_3 \in \mathbb{R}^2$ ,  $H_{3,j} \in \mathbb{R}^{1 \times 2}$ , and  $w_{3,j} \in \mathbb{R}^2$  is white Gaussian noise with zero mean and covariance matrix  $I_2$ .

Therefore, the state-space model of (42) can be obtained by combining the state-space models (A.18), (A.23) and (A.24) as follows

$$\begin{bmatrix} s_{1,j+1} \\ s_{2,j+1} \\ s_{3,j+1} \end{bmatrix} = \begin{bmatrix} F_1 & 0 & 0 \\ 0 & F_2 & 0 \\ 0 & 0 & F_3 \end{bmatrix} \begin{bmatrix} s_{1,j} \\ s_{2,j} \\ s_{3,j} \end{bmatrix} + \begin{bmatrix} G_1 & 0 & 0 \\ 0 & G_2 & 0 \\ 0 & 0 & G_3 \end{bmatrix} \begin{bmatrix} w_{1,j} \\ w_{2,j} \\ w_{3,j} \end{bmatrix},$$

$$y_j = \begin{bmatrix} H_{1,j} & H_{2,j} & H_{3,j} \end{bmatrix} \begin{bmatrix} s_{1,j}^T & s_{2,j}^T & s_{3,j}^T \end{bmatrix}^T.$$

## References

- [1] B. D. O. Anderson and J. B. Moore. *Optimal Filtering*. Prentice Hall, New Jersey, 1979.
- [2] K. J. Åström. *Introduction to stochastic control theory*. Courier Corporation, 2012.
- [3] G. Atluri, A. Karpatne, and V. Kumar. Spatio-temporal data mining: A survey of problems and methods. *ACM Computing Surveys (CSUR)*, 51(4):1–41, 2018.
- [4] J. Bergstra, R. Bardenet, Y. Bengio, and B. Kégl. Algorithms for hyper-parameter optimization. In J. Shawe-Taylor, R. Zemel, P. Bartlett, F. Pereira, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems*, volume 24. Curran Associates, Inc., 2011.
- [5] G. Bottegal and G. Pillonetto. The generalized cross validation filter. *Automatica*, 90:130–137, 2018.
- [6] J. V. Candy. *Model-based signal processing*, volume 36. John Wiley & Sons, 2005.
- [7] A. Carron, M. Todescato, R. Carli, L. Schenato, and G. Pillonetto. Machine learning meets Kalman filtering. In *2016 IEEE 55th Conference on Decision and Control (CDC)*, pages 4594–4599. IEEE, 2016.
- [8] C. Chen. *Linear system theory and design*. Oxford University Press, New York, 3 edition, 1999.
- [9] T. Chen. On kernel design for regularized LTI system identification. *Automatica*, 90:109–122, 2018.
- [10] T. Chen and M. S. Andersen. On semiseparable kernels and efficient implementation for regularized system identification and function estimation. *Automatica*, 132:109682, 2021.

- [11] T. Chen and L. Ljung. Implementation of algorithms for tuning parameters in regularized least squares problems in system identification. *Automatica*, 49(7):2213–2220, 2013.
- [12] T. Chen, H. Ohlsson, and L. Ljung. On the estimation of transfer functions, regularizations and gaussian processes—revisited. *Automatica*, 48(8):1525–1535, 2012.
- [13] A. E. Gelfand, P. Diggle, P. Guttorp, and M. Fuentes. *Handbook of spatial statistics*. CRC press, 2010.
- [14] M. G. Genton. Classes of kernels for machine learning: a statistics perspective. *Journal of machine learning research*, 2(Dec):299–312, 2001.
- [15] T. Glad and L. Ljung. *Control theory: Multivariable and nonlinear methods*. Taylor & Francis, 2000.
- [16] T. Hofmann, B. Schölkopf, and A. J. Smola. Kernel methods in machine learning. *The annals of statistics*, 36(3):1171–1220, 2008.
- [17] Y. Kuang, T. Chen, F. Yin, and R. Zhong. Recursive implementation of Gaussian process regression for spatial-temporal data modeling. In *2019 11th International Conference on Wireless Communications and Signal Processing (WCSP)*, pages 1–7. IEEE, 2019.
- [18] Q. Liu, H. S. Abbas, and J. M. Velni. An LMI-based approach to distributed model predictive control design for spatially-interconnected systems. *Automatica*, 95:481–487, 2018.
- [19] L. Ljung. *System Identification - Theory for the User*. Prentice-Hall, Upper Saddle River, N.J., 2nd edition, 1999.
- [20] L. Ljung. *System Identification Toolbox for use with MATLAB. Version 5*. The MathWorks, Inc, Natick, MA, 5th edition, 2000.
- [21] M. J. Menne, I. Durre, B. Korzeniewski, S. McNeal, K. Thomas, X. Yin, S. Anthony, R. Ray, R. S. Vose, B. E. Gleason, et al. Global historical climatology network-daily (GHCN-daily), version 3. *NOAA National Climatic Data Center*, 10:V5D21VHZ, 2012.
- [22] N. Pelekis, B. Theodoulidis, I. Kopanakis, and Y. Theodoridis. Literature review of spatio-temporal database models. *The Knowledge Engineering Review*, 19(3):235–274, 2004.
- [23] F. Perez-Cruz, S. Van Vaerenbergh, J. J. Murillo-Fuentes, M. Lazaro-Gredilla, and I. Santamaria. Gaussian processes for nonlinear signal processing: An overview of recent advances. *IEEE Signal Processing Magazine*, 30(4):40–50, 2013.
- [24] G. Pillonetto, F. Dinuzzo, T. Chen, G. De Nicolao, and L. Ljung. Kernel methods in system identification, machine learning and function estimation: A survey. *Automatica*, 50(3):657–682, 2014.
- [25] J. Quinonero-Candela and C. E. Rasmussen. A unifying view of sparse approximate Gaussian process regression. *The Journal of Machine Learning Research*, 6:1939–1959, 2005.
- [26] S. Särkkä, A. Solin, and J. Hartikainen. Spatio-temporal learning via infinite-dimensional Bayesian filtering and smoothing. *IEEE Signal Processing Magazine*, 30(4):51–61, 2013.
- [27] M. Todescato, A. Carron, R. Carli, G. Pillonetto, and L. Schenato. Efficient spatio-temporal Gaussian regression via Kalman filtering. *Automatica*, 118:109032, 2020.
- [28] N. Wahlström, P. Axelsson, and F. Gustafsson. Discretizing stochastic dynamical systems using Lyapunov equations. *IFAC Proceedings Volumes*, 47(3):3726–3731, 2014.
- [29] C. K. Williams and C. E. Rasmussen. *Gaussian processes for machine learning*, volume 2. MIT press Cambridge, MA, 2006.
- [30] A. Wilson and R. Adams. Gaussian process kernels for pattern discovery and extrapolation. In *International conference on machine learning*, pages 1067–1075. PMLR, 2013.
- [31] J. Zhang, Y. Kuang, T. Chen, X. Lu, F. Yin, and R. Zhong. Efficient recursive implementation of spatial-temporal Gaussian process regression. In *2020 39th Chinese Control Conference (CCC)*, pages 1081–1086. IEEE, 2020.
- [32] M. Zorzi. Autoregressive identification of kronecker graphical models. *Automatica*, 119:109053, 2020.