

Compressible Spectral Mixture Kernels with Sparse Dependency Structures for Gaussian Processes

Kai Chen^a, Feng Yin^b (✉), and Shuguang Cui^{b,c}

^a*School of Mathematics and Statistics, Central South University, Changsha 410083, China*

^b*School of Science and Engineering (SSE), The Chinese University of Hong Kong, Shenzhen 518172, China*

^c*Future Network of Intelligence Institute (FNii), The Chinese University of Hong Kong, Shenzhen 518172, China*

Abstract

Spectral mixture (SM) kernels comprise a powerful class of generalized kernels for Gaussian processes (GPs) to describe complex patterns. This paper introduces model compression and time- and phase (TP) modulated dependency structures to the original (SM) kernel for improved generalization of GPs. Specifically, by adopting Bienaymés identity, we generalize the dependency structure through cross-covariance between the SM components. Then, we propose a novel SM kernel with a dependency structure (SMD) by using cross-convolution between the SM components. Furthermore, we ameliorate the expressiveness of the dependency structure by parameterizing it with time and phase delays. The dependency structure has clear interpretations in terms of spectral density, covariance behavior, and sampling path. To enrich the SMD with effective hyperparameter initialization, compressible SM kernel components, and sparse dependency structures, we introduce a novel structure adaptation (SA) algorithm in the end. A thorough comparative analysis of the SMD on both synthetic and real-life applications corroborates its efficacy.

Keywords: Gaussian processes, spectral mixture, dependency structure, time and phase delays, structure adaptation

1. Introduction

Gaussian processes (GPs) constitute an important class of Bayesian nonparametric models for machine learning [1]. A GP models an underlying system by applying a Gaussian prior to the underlying function and computes the posterior distribution over this function given the observations. This allows GPs to learn a function approximation well if a sufficient number of observations is accumulated. Furthermore, GPs can avoid overfitting in cases where only a little evidence [2, 3] is available. A GP can model a large class of systems by selecting and designing the kernel function, which reflects the autocovariance structure of the system. However, similar to other kernel learning methods, such as support

vector machines (SVMs), selecting kernel function is one of the most important factors for GP model because an expressive kernel determines the representation ability of GP, and the posterior distribution can change significantly by using different kernels.

For GPs, however, a kernel is usually selected subjectively, heavily depending on expert knowledge and empirical analysis of data patterns. There are a handful of base kernels and their combinations for diversified GP learning applications, such as received signal strength (RSS)-based radio map modeling [4] and wireless traffic prediction [5, 6, 7]. To avoid human intervention, automatic [8, 9, 10] and generalized kernel designs [11] are highly demanded for GPs. Since [8] introduced an automatic and expressive kernel, called spectral mixture (SM) kernel, by modeling the spectral density (SD) of a stationary signal with a sum of Gaussians. GPs with the SM kernel have been successfully employed in various applications, such as medical time series prediction [12], Arctic coastal erosion forecasting [13], and urban environmental monitoring using sensor networks [14].

Briefly, there are many advances focusing on diversified extensions of SM kernel rather than the latent dependency structure between SM components. The SM product (SMP) kernel [15] constructed via a product of several SM kernels on the input dimensions can model image and spatial data. The non-stationary SM (NSM) kernel [16, 17] demonstrates a compelling ability to represent input-dependent covariances between inputs. The grid SM (GSM) kernel [10] linearly combines independent low-rank sub-kernels to approximate the underlying covariance of a stationary signal. However, all these variants of SM kernel cannot explicitly represent the dependency structures.

In this paper, we show that there are extensive dependency structures between the SM components. We propose a new SM kernel with a dependency structure (SMD). In addition, we introduce a structure adaptation (SA) algorithm to compress the SMD effectively. We demonstrate the benefits of modeling the dependency structure and applying model compression.

More specifically, we generalize the dependency structure as a cross covariance between the SM components by using Bienaymé’s identity for the linear superposition form of GP. The cross covariance is constructed by a cross convolution between the SM components. To ameliorate the expressiveness of the dependency structure, we design a complex-valued Gaussian SD that incorporates both time and phase (TP) delays for SM component. We then construct a positive definite SMD kernel that handles the dependency structures. The spectral density, covariance behavior, and sampling path of the dependency structure are interpretable and informative. We propose a SA algorithm, including bootstrap-based hyperparameter initialization (BHI), compression of SMD components, and sparsification of the dependency structures, to improve the learning efficiency and interpretability of the SMD automatically. The SA algorithm can prevent the hyperparameter space from expansion and retain valuable dependency structures of SMD. To analyze the dependency structure, we introduce a measure of dependency intensity $\gamma_{i,j}$ (see Eq. (10)). The SMD kernel can be regarded as the generalization of the original SM kernel; that is, by only considering the autocovariance of its components, the SMD kernel will

reduce to the SM kernel. Preliminary results of this work have been presented in [18]. While this paper includes abundant new contributions, such as:

- We represent the dependency structure as a generalized cross covariance.
- We construct a complex-valued Gaussian mixture model (GMM) characterized by TP delays to model the SDs.
- A new SMD kernel demonstrates good interpretability, and more importantly better expressiveness.
- An effective and interpretable SA algorithm for compressible SMD.
- We investigate the interpolation and extrapolation performances, scalable learning, dependency structures, compression ratio (CR), and sparsity ratio (SR), of the SMD on multiple synthetic and real-life datasets.

The rest of this paper is organized as follows. Background on GPs and SM kernels and a summary of the related works are given in Section 2. In Section 3, we present our motivation. Section 4 introduces our SMD kernel. Section 5 describes the SA algorithm for the SMD. Section 6 shows multiple experiments. Concluding remarks and future works are discussed in Section 7.

2. Background and related work

In this section, we review GPs, SM kernels, and related works.

2.1. Gaussian processes

Given an observation pair $\{\mathbf{x}_i \in \mathbb{R}^P, y_i\}$ with P dimensional input, a GP can represent a function map $y = f(\mathbf{x}) + \epsilon$ to approximate the underlying system, where ϵ is the noise. In this paper, we mainly consider real-valued GPs. From the function-space view, a GP [1, 3] defines a prior distribution over functions, completely specified by its first-order and second-order statistics, namely, the mean function $m(\mathbf{x})$ and the covariance function $k(\mathbf{x}, \mathbf{x}')$. In general, a GP is defined as $f(\mathbf{x}) \sim \mathcal{GP}(m(\mathbf{x}), k(\mathbf{x}, \mathbf{x}'))$. We use the terms covariance function, kernel, and kernel function interchangeably.

Given a selected kernel function $k(\mathbf{x}, \mathbf{x}')$ and training set $\mathcal{D} = \{X, \mathbf{y}\}$ for a GP model, we can analytically compute the predictive mean $\tilde{\mathbf{y}}^*$ and variances $\mathbb{V}[\mathbf{y}^*]$ (that is, its predictive uncertainty) for the testing set X^* by using the inferred posterior distribution $p(\mathbf{y}^*|X^*, X, \mathbf{y}) \sim \mathcal{N}(\tilde{\mathbf{y}}^*, \mathbb{V}[\mathbf{y}^*])$. $k(\mathbf{x}, \mathbf{x}')$ has free hyperparameters Θ determining the model complexity of a GP. A GP model is optimized by minimizing the negative log-marginal likelihood (NLML), $\mathcal{L}_{\text{NLML}} \triangleq -\log p(\mathbf{y}|X, \Theta)$, where \mathcal{L} is obtained through marginalization over the latent function $f(X)$ [1, 3].

2.2. Spectral mixture kernels

The SM kernel [8] has been derived with the aid of Bochner’s Theorem [19, 20]. The essence of this theorem is that a function k on \mathbb{R}^P is the covariance function of a weakly stationary mean square continuous complex-valued random process on \mathbb{R}^P if and only if it can be represented as $k(\boldsymbol{\tau}) = \int_{\mathbb{R}^P} e^{2\pi i \mathbf{s}^\top \boldsymbol{\tau}} d\psi(\mathbf{s})$, where ψ is a positive finite measure and i denotes the imaginary unit. If ψ has density $\hat{k}(\mathbf{s})$, then \hat{k} is called the SD or power spectrum of the kernel. Moreover, k and \hat{k} constitute an Fourier transform (FT) pair; that is, $\hat{k}(\mathbf{s}) = \mathcal{F}_{\boldsymbol{\tau} \rightarrow \mathbf{s}}[k(\boldsymbol{\tau})](\mathbf{s})$ and $k(\boldsymbol{\tau}) = \mathcal{F}_{\mathbf{s} \rightarrow \boldsymbol{\tau}}^{-1}[\hat{k}(\mathbf{s})](\boldsymbol{\tau})$, where the operators $\mathcal{F}_{\boldsymbol{\tau} \rightarrow \mathbf{s}}$ and $\mathcal{F}_{\mathbf{s} \rightarrow \boldsymbol{\tau}}^{-1}$ denote the FT and the inverse FT, respectively. Originally, the SM kernel [8] was constructed by approximating the underlying SD as a mixture of Q Gaussians in the frequency domain. The SM kernel can approximate any stationary kernel with a sufficient number of Gaussian components. By applying the inverse FT, we can obtain the SM kernel as $k_{\text{SM}} = \sum_{i=1}^Q w_i k_{\text{SM},i}(\boldsymbol{\tau})$ with

$$k_{\text{SM},i}(\boldsymbol{\tau}) = \exp(-2\pi^2 \boldsymbol{\tau}^\top \Sigma_i \boldsymbol{\tau}) \cos(2\pi \boldsymbol{\tau}^\top \boldsymbol{\mu}_i), \quad (1)$$

where $k_{\text{SM},i}(\boldsymbol{\tau}) = \mathcal{F}_{\mathbf{s} \rightarrow \boldsymbol{\tau}}^{-1}[\hat{k}_{\text{SM},i}(\mathbf{s})](\boldsymbol{\tau})$ is the i -th SM component, w_i , $\boldsymbol{\mu}_i$, and Σ_i are the signal magnitude, center frequency, and frequency bandwidth parameters of $k_{\text{SM},i}(\boldsymbol{\tau})$, respectively.

2.3. Related work

There is rich literature on GPs related to SM kernel function design and analysis [3, 8, 9, 21, 22]. This section mainly focuses on the family of SM kernels [8, 15, 10] and some new variants. More related SM extensions are surveyed in [16, 18, 10]. Similar to the compositional form of the SM kernel, additive GPs [9, 21] implicitly sum over some one-dimensional base kernels to construct a flexible kernel representation. For the advances of SM kernel [15, 23], we have the SMP kernel with $k_{\text{SMP}}(\boldsymbol{\tau}|\Theta) = \sum_{i=1}^Q \prod_{p=1}^P k_{\text{SM}}(\tau_p|\Theta_p)$ and NSM kernel [16, 17]. The NSM kernel includes a non-stationary Gibbs kernel $k_{\text{Gibbs},i}(x, x')$ replacing the exponential part of the SM kernel, and input-dependent $w_i(x)$ and $\mu_i(x)$ corresponding to w_i and μ_i in the SM kernel. Recently, an approach was proposed in [18] that encodes simple dependency structures between components of an SM kernel. However, similar to the existing additive GPs, most of the existing SM variants assume that GP components specified by SM components are independent and ignore their possible dependency structure.

On the other hand, a quantification of the dependency structure between components in GPs was initially proposed in [21]. However, no further investigation in modeling the dependency structure is presented therein. In [18], the main challenges in sparsifying the dependency structure and modeling TP delays of dependency structure are still unsolved.

In [10], another hyperparameter efficient inference approach fixes $\boldsymbol{\mu}_i$ and Σ_i of SM and allows only w_i to be optimizable, which could result in a sparse SM kernel. In [24], a Lévy process was introduced to automatically select the number of SM components, but the selection is not stable and easily encounters

overfitting. In short, the inference and optimization of the SMD in terms of hyperparameter initialization, model compression, and dependency structure sparsity have not yet been studied.

3. Motivation

This section aims to give a generalized SM kernel with dependency structure and further comment on its properties. All components are additive for the original SM kernel [8]. Any function f drawn from a GP with the SM kernel k_{SM} , that is, $f \sim \mathcal{GP}(0, k_{\text{SM}})$, can be described as $f = \sum_{i=1}^Q f_{\text{SM},i}$, where $f_{\text{SM},i} \sim \mathcal{GP}(0, w_i k_{\text{SM},i})$. To simplify our notations, we use $\mathbf{f}_{\text{SM},i}$ and $\mathbf{f}_{\text{SM},i}^*$ to denote the respective function values evaluated on X and X^* , respectively.

Generalization of dependency structure: Generally, by using the Bienaymé’s identity [25, 26] for the linear form of f , the generalized covariance of f is given by

$$\begin{aligned} \mathbb{V}[f] &= \sum_{i=1}^Q \sum_{j=1}^Q \text{Cov}(f_{\text{SM},i}, f_{\text{SM},j}) \\ &= \underbrace{\sum_{i=j} \mathbb{V}[f_{\text{SM},i}]}_{\text{autocovariance}} + \underbrace{\sum_{i \neq j} \text{Cov}(f_{\text{SM},i}, f_{\text{SM},j})}_{\text{cross covariance}}, \end{aligned} \quad (2)$$

where $\mathbb{V}[f_{\text{SM},i}] = \text{Cov}(f_{\text{SM},i}, f_{\text{SM},i})$ is the autocovariance. Here, the autocovariance of random function $f_{\text{SM},i}$ is computed as $\text{Cov}(f_{\text{SM},i}, f_{\text{SM},i}) = w_i k_{\text{SM},i}$. Therefore, we reformulate Eq. (2) as

$$\mathbb{V}[f] = \sum_{i=1}^Q w_i k_{\text{SM},i} + \overbrace{\sum_{i=1}^Q \sum_{j \neq i}^Q \text{Cov}(f_{\text{SM},i}, f_{\text{SM},j})}^{\text{dependency structure}}. \quad (3)$$

For the SM kernel, however, it restricts that $f_{\text{SM},i}$ and $f_{\text{SM},j}$ are independent with $\text{Cov}(f_{\text{SM},i}, f_{\text{SM},j}) = 0, \forall (i \neq j)$. Unfortunately, there is no evident support that $\text{Cov}(f_{\text{SM},i}, f_{\text{SM},j}) = 0$. In this paper, we consider $\text{Cov}(f_{\text{SM},i}, f_{\text{SM},j}) \neq 0$ as a general dependency structure to free the SM kernel.

4. Spectral mixture kernel with dependency structure

In this section, we propose for the first time an extended SM kernel incorporating dependency structure and its TP augmentations.

4.1. Modeling dependency structure using convolution

A stationary covariance function $k(\mathbf{x}, \mathbf{x}')$ can be represented in convolution form on \mathbb{R}^P , as in [27, 28], $k(\mathbf{x}, \mathbf{x}') \triangleq \int_{\mathbb{R}^P} g(\mathbf{u}) g(\boldsymbol{\tau} - \mathbf{u}) d\mathbf{u} = (g * g)(\boldsymbol{\tau})$, where $\boldsymbol{\tau} \triangleq \mathbf{x} - \mathbf{x}'$ and $*$ denotes the convolution operator. Since convolution in the time domain corresponds to multiplication in the frequency domain, we have the squared form of the i -th SM component as

$$\begin{aligned} w_i \hat{k}_{\text{SM},i}(\mathbf{s}) &= \mathcal{F}_{\boldsymbol{\tau} \rightarrow \mathbf{s}}[(g_{\text{SM},i} * g_{\text{SM},i})(\boldsymbol{\tau})](\mathbf{s}) \\ &= \hat{g}_{\text{SM},i}^2(\mathbf{s}), \end{aligned} \quad (4)$$

where $\hat{g}_{\text{SM},i}(\mathbf{s})$ is the SD of the i -th SM basis component. For the dependency structures $\text{Cov}(f_{\text{SM},i}, f_{\text{SM},j})$, one possible approach is to employ the cross correlation between functions $f_{\text{SM},i}$ and $f_{\text{SM},j}$, which is equal to the convolution of the two (weighted) kernels $w_i k_{\text{SM},i}$ and $w_j k_{\text{SM},j}$, namely,

$$f_{\text{SM},i} \star f_{\text{SM},j} = \mathcal{F}_{\mathbf{s} \rightarrow \boldsymbol{\tau}}^{-1} [w_i \varphi_{\text{SM},i}(\mathbf{s}) \cdot \overline{w_j \varphi_{\text{SM},j}(\mathbf{s})}](\boldsymbol{\tau}), \quad (5)$$

where $\varphi_{\text{SM},i}(\mathbf{s}) = \mathcal{N}(\mathbf{s}; \boldsymbol{\mu}_i, \Sigma_i)$, $\mathcal{F}_{\mathbf{s} \rightarrow \boldsymbol{\tau}}^{-1}$, \star , and $\overline{(\cdot)}$ denote the inverse FT, the cross-correlation operator, and the complex conjugate operator, respectively. However, by directly using Eq. (5), we will obtain a kernel that is the inverse FT of the squared Gaussian $w_i^2 \varphi_{\text{SM},i}^2(\mathbf{s})$ when $i = j$. This is different from the original SM component. To ensure the compatibility between the dependency structure and SM component when $i = j$, we, therefore, consider the convolution between the *basis components* $g_{\text{SM},i}(\boldsymbol{\tau})$ and $g_{\text{SM},j}(\boldsymbol{\tau})$, $g_{\text{SM},i}(\boldsymbol{\tau}) * g_{\text{SM},j}(\boldsymbol{\tau})$, where $g_{\text{SM},i}(\boldsymbol{\tau}) = \mathcal{F}_{\mathbf{s} \rightarrow \boldsymbol{\tau}}^{-1}[\hat{g}_{\text{SM},i}(\mathbf{s})](\boldsymbol{\tau})$. Thus, we can describe $\text{Cov}(f_{\text{SM},i}, f_{\text{SM},j})$ well. Note that $\text{Cov}(f_{\text{SM},i}, f_{\text{SM},j})$ does not introduce additional parameters for the dependency structure.

4.2. Time and phase characterized Gaussian spectral density

In signal processing, for a signal sample of an underlying process, time delay differences between different signal frequency components characterize their temporal relationship and influence the signal's shape. Furthermore, due to the nature of FT, a signal in the frequency domain always has a complex representation with magnitude (real) and phase (imaginary) parts. It can be of interest to know not only the magnitude but also the phase of the SD. The phase difference is useful for understanding the interference (dependency structure) phenomenon between components of a physical process. However, the dependency structure investigated in [18] only paints a picture of the magnitude difference between the SM components. Here, we introduce TP parameterization for the SD of the SM component to enrich its representation capacity.

Based on the property of FT, shifting a signal $k(\boldsymbol{\tau})$ with time delay $\boldsymbol{\theta}$ in the time domain is equivalent to multiplying a complex exponential in the frequency domain i.e., $\hat{k}_{\boldsymbol{\theta}}(\mathbf{s}) = e^{-2\pi \boldsymbol{\theta} \mathbf{s}^T} \hat{k}(\mathbf{s})$, where $k_{\boldsymbol{\theta}}(\boldsymbol{\tau}) \triangleq k(\boldsymbol{\tau} - \boldsymbol{\theta})$ and $\hat{k}(\mathbf{s}) = \mathcal{F}_{\boldsymbol{\tau} \rightarrow \mathbf{s}}[k(\boldsymbol{\tau})](\mathbf{s})$ [29]. For any phase delay function $k_{\boldsymbol{\phi}}(\boldsymbol{\tau})$ with a phase delay vector $\boldsymbol{\phi}$, the FT of $k_{\boldsymbol{\phi}}(\boldsymbol{\tau})$ in the frequency domain is $\hat{k}_{\boldsymbol{\phi}}(\mathbf{s}) = e^{-2\pi \boldsymbol{\phi}^T \mathbf{s}} \hat{k}(\mathbf{s})$. For

the SMD kernel, we can directly embed the time delay $\boldsymbol{\theta}_i$ and phase delay $\boldsymbol{\phi}_i$ into the SD, $\hat{k}_{\text{SM},i}(\mathbf{s})$, yielding the following complex-valued *TP delay SD* function:

$$\hat{k}_{\text{SMD},i}(\mathbf{s}) = w_i \varphi_{\text{SM},i}(\mathbf{s}) \overbrace{\exp(-2\pi\iota(\boldsymbol{\theta}_i \mathbf{s} + \boldsymbol{\phi}_i))}^{\text{TP delays}}. \quad (6)$$

4.3. Time- and phase modulated dependency structure

Considering the TP modulated SD function in Eq. (6) and adopting the squared form in Eq. (4), we define $\hat{g}_{\text{SMD},i}(\mathbf{s}) = \hat{k}_{\text{SMD},i}^{1/2}(\mathbf{s})$. We then express the corresponding SD with the dependency structure as

$$\begin{aligned} \hat{k}_{\text{SMD}}^{i \times j}(\mathbf{s}) &= \hat{g}_{\text{SMD},i}(\mathbf{s}) \cdot \overline{\hat{g}_{\text{SMD},j}(\mathbf{s})} \\ &= w_{ij} a_{ij} \varphi_{\text{SMD},ij}(\mathbf{s}) \overbrace{\exp(-\pi\iota(\boldsymbol{\theta}_{ij} \mathbf{s} + \boldsymbol{\phi}_{ij}))}^{\text{cross TP delays}}, \end{aligned} \quad (7)$$

with the following parameters:

- cross weight: $w_{ij} = \sqrt{w_i w_j}$,
- cross amplitude: $a_{ij} = |4\pi^2 \Sigma_i \Sigma_j|^{-\frac{1}{4}} \mathcal{N}(\boldsymbol{\mu}_i; \boldsymbol{\mu}_j, \frac{\Sigma_i + \Sigma_j}{2})$,
- cross Gaussian: $\varphi_{\text{SMD},ij}(\mathbf{s}) = \mathcal{N}(\mathbf{s}; \boldsymbol{\mu}_{ij}, \Sigma_{ij})$
- cross mean of $\varphi_{\text{SMD},ij}$: $\boldsymbol{\mu}_{ij} = \frac{\Sigma_i \boldsymbol{\mu}_j + \Sigma_j \boldsymbol{\mu}_i}{\Sigma_i + \Sigma_j}$,
- cross covariance of $\varphi_{\text{SMD},ij}$: $\Sigma_{ij} = \frac{2\Sigma_i \Sigma_j}{\Sigma_i + \Sigma_j}$,
- cross time delay: $\boldsymbol{\theta}_{ij} = \boldsymbol{\theta}_i - \boldsymbol{\theta}_j$,
- cross phase delay: $\boldsymbol{\phi}_{ij} = \boldsymbol{\phi}_i - \boldsymbol{\phi}_j$.

The cross amplitude a_{ij} is a normalization constant that only depends on the difference between components i and j . Without TP delays, the SD term in Eq. (7) can be reduced as follows: $\hat{k}_{\text{SMD},\boldsymbol{\theta}=0,\boldsymbol{\phi}=0}^{i \times j}(\mathbf{s}) \triangleq \hat{g}_{\text{SM},i}(\mathbf{s}) \cdot \overline{\hat{g}_{\text{SM},j}(\mathbf{s})}$ [18].

Remark 1. According to Eq. (7), the closer the components are, the larger the weight w_{ij} , frequency $\boldsymbol{\mu}_{ij}$, and scale Σ_{ij} , and the greater the dependency structure in the SMD.

4.4. Spectral mixture with TP modulated dependency structure

In light of the SM kernel, by applying the inverse FT, we can define the dependency structure as

$$\begin{aligned} k_{\text{SMD}}^{i \times j} &= \mathcal{F}_{\mathbf{s} \rightarrow \boldsymbol{\tau}}^{-1} [\hat{g}_{\text{SMD},i}(\mathbf{s}) \cdot \overline{\hat{g}_{\text{SMD},j}(\mathbf{s})}] (\boldsymbol{\tau}) \\ &= c_{ij} \exp\left(-\frac{\boldsymbol{\tau}_\theta^\top \Sigma_{ij} \boldsymbol{\tau}_\theta}{2}\right) \exp(\iota(\boldsymbol{\tau}_\theta^\top \boldsymbol{\mu}_{ij} - \boldsymbol{\phi}_{ij} \pi)), \end{aligned} \quad (8)$$

where $\tau_\theta \triangleq 2\pi(\tau - \frac{\theta_{ij}}{2})$ is the Euclidean distance with time delay and $c_{ij} = w_{ij}a_{ij}$ is the normalization term incorporating the cross weight and cross amplitude and it does not depend on τ . Note that c_{ij} indicates the largest degree of the dependency structure because the exponential term has a max value of 1.

Given an SM kernel with Q components, we can obtain the corresponding dependency structures by considering the symmetric properties of SD as follows:

$$k_{\text{SMD}} = \sum_{i=1}^Q \sum_{j=1}^Q c_{ij} \exp\left(-\frac{\tau_\theta^\top \Sigma_{ij} \tau_\theta}{2}\right) \cos(\tau_\theta^\top \boldsymbol{\mu}_{ij} - \phi_{ij}\pi). \quad (9)$$

The positive semidefinite (PSD) property of the SMD kernel is equivalent to saying that its SD, $\hat{k}_{\text{SMD}}(\mathbf{s})$, is PSD as well [19, 20]. Given any finite set of non-zero vectors $[\mathbf{z}_1, \dots, \mathbf{z}_N]^\top \in \mathbb{C}^{N \times P}$ with complex entry, $\mathbf{s} \in \mathbb{R}^P$, we have $\sum_{n=1}^N \left| \sum_{i=1}^Q \mathbf{z}_n \hat{g}_{\text{SMD},i}(\mathbf{s}) \right|^2 \geq 0$. Hence, the SMD kernel must be PSD. We have $\text{Cov}(f_{\text{SM},i}, f_{\text{SM},j}) = k_{\text{SMD}}^{i \times j}(\boldsymbol{\tau})$ to represent the dependency structure in Eq. (3). There are Q^2 structures with Q original components plus $Q^2 - Q$ dependency structures.

Quantification of the dependency structure: To measure the intensity of the dependency structure, we normalize the dependency structure as

$$\gamma_{ij}(\boldsymbol{\tau}) = \frac{k_{\text{SMD}}^{i \times j}(\boldsymbol{\tau})}{\sqrt{w_i k_{\text{SM},i}(\boldsymbol{\tau}) \cdot w_j k_{\text{SM},j}(\boldsymbol{\tau})}}. \quad (10)$$

Note that γ_{ij} has a range with $[-1, 1]$. For $i = j$, we have $\gamma_{ij} = 1$ when $w_i k_{\text{SM},i}(\boldsymbol{\tau}) > 0$ and $\gamma_{ij} = -1$ when $w_i k_{\text{SM},i}(\boldsymbol{\tau}) < 0$.

4.5. Interpretation of dependency structure

In Fig. 1, we show the covariances, SDs, sampling paths, and posterior distributions in terms of amplitude, peak, and trend between the SM (dashed red) and SMD (dashed blue) kernel. The differences between SM and SMD are clear. Without TP delays, subplots (b) and (g) of Fig. 1 show that the dependency structure can reinforce the magnitudes of both SD and covariance in SM (shown in subplots (a) and (f)) but does not change the decaying behavior of covariance a lot. The dependency structure in the frequency domain (subplot (g)) is the intersection (modeled as a Gaussian, see Eq. (7)) between two SM components. When having TP delays, the dependency structure can reinforce or weaken the covariances and SDs of the original kernel. In subplots (c) and (e), the covariance range of SMD is much extended and larger than SM due to the time delay. Specifically, the dependency structure can largely change the magnitudes of both SD (shown in subplots (h), (i), and (j)) and covariance (shown in subplots (c), (d), and (e)), shapes of SD, and decaying behaviors of the covariance, and further reduce the predictive uncertainties (show in subplots (m), (n), and (o)).

Given six observations (marked with black crosses) and conditions on them, the learned posterior distribution and sampling path are shown in subplots

(k), (l), (m), (n), and (o). Interestingly, due to the dependency structure, the predictive confidence interval (CI) of the SMD is significantly tighter (in blue shadow) than that of the SM (in red shadow).

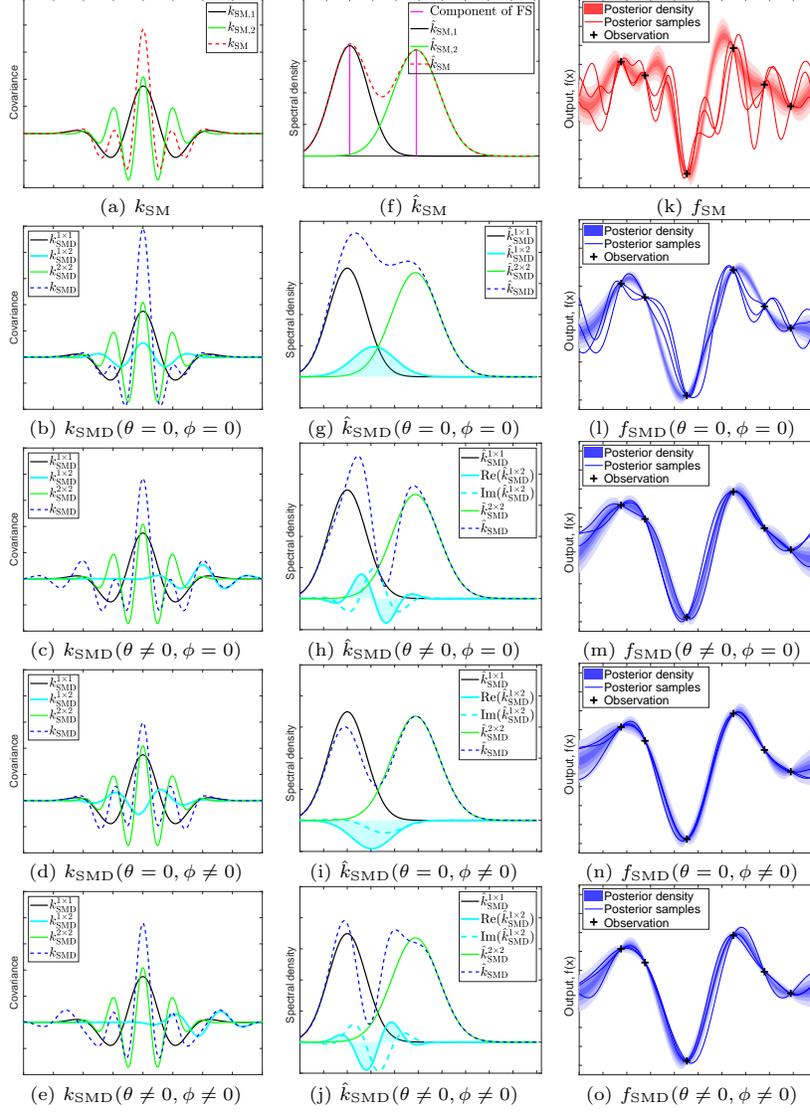


Figure 1: Covariances (the first column), SDs (the second column), sampling path (the third column), and posterior distributions based on GPs with the SM ($Q = 2$) and SMD ($Q = 2$) kernels conditioned on six observations. The samples of all GP models were obtained using 200 equally spaced points.

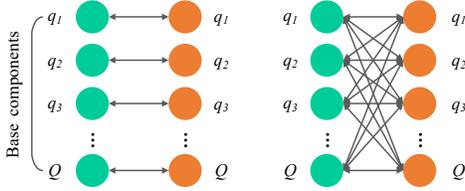


Figure 2: SM kernel (left) vs SMD kernel (right) with Q original components, where $q_i=\{1,\dots,Q\}$ denotes the i -th SM component. The SM models only the autocovariance between the component itself. The SMD models both auto- and cross-covariance between components.

4.6. Comparisons between the SMD and related kernels

In Fig. 2, we visualize the covariance differences between the SM and SMD (shown in Eq. (3)), where each link (in black solid) represents a covariance structure of the kernel. Circle q_i corresponds to a $f_{\text{SM},i}$. The cross connection denotes a $\text{Cov}(f_{\text{SM},i}, f_{\text{SM},j})$ of $f_{\text{SM},i}$ and $f_{\text{SM},j}$. The SM considers only the autocovariance $\text{Cov}(f_{\text{SM},i}, f_{\text{SM},i})$ of its components and ignores their dependency structures. Table 1 summarizes the differences between the SMD and SM kernels in terms of their hyperparameters for a P -dimensional input setting. For NSM, each hyperparameter of the original SM is parameterized as a GP with a squared exponential (SE) kernel, for instance, the weight w_i becomes $w_{i,\mathbf{x}} \sim \mathcal{GP}(0, k_{\text{SE}}(\mathbf{x}, \mathbf{x}'))$ in NSM. Thus NSM needs three times more hyperparameters than SM because $w_{i,\mathbf{x}}$ usually has three hyperparameters. Without TP delays, the hyperparameter space of the SMD is equal to that of the SM. The price paid for incorporating TP delays in the SMD is that the gradient computation is more involved because of additional TP hyperparameters.

Table 1: Comparisons between the SMD and other SM kernels in terms of hyperparameters and the number of hyperparameters. For an initial large Q in the SM and SMD, the number of components retained after compression, Q_{rest} , is much smaller than Q .

Kernel	hyperparameters	Number of hyperparameters
SM	$\{w_i, \boldsymbol{\mu}_i, \Sigma_i\}_{i=1}^Q$	$(2P+1)Q$
NSM	$\{w_{i,\mathbf{x}}, \boldsymbol{\mu}_{i,\mathbf{x}}, \Sigma_{i,\mathbf{x}}\}_{i=1}^Q$	$3 \times (2P+1)Q$
SMD $\phi=0, \boldsymbol{\theta}=0$	$\{w_i, \boldsymbol{\mu}_i, \Sigma_i\}_{i=1}^Q$	$(2P+1)Q$
SMD	$\{w_i, \boldsymbol{\mu}_i, \Sigma_i, \boldsymbol{\theta}_i, \boldsymbol{\phi}_i\}_{i=1}^Q$	$(4P+1)Q$
Compressed SMD	$\{w_i, \boldsymbol{\mu}_i, \Sigma_i, \boldsymbol{\theta}_i, \boldsymbol{\phi}_i\}_{i=1}^{Q_{\text{rest}}}$	$(4P+1)Q_{\text{rest}}$
SMD with SA	$\{w_i, \boldsymbol{\mu}_i, \Sigma_i, \boldsymbol{\theta}_i, \boldsymbol{\phi}_i\}_{i=1}^{Q_{\text{rest}}}$	$(1+2P(2-\alpha_{\text{SR}}))Q_{\text{rest}}$

5. Structure adaptation for the spectral mixture with dependency structure

The SM kernel has been known for its large number of hyperparameters (with size $3Q$) [8]. This complicates the inference, learning, and interpretability

Algorithm 1: Structure adaption for SMD

Input : Initial number of components, Q_{init} , number of training attempts, M_{init} , number of pruned components, $Q_{\text{prune}} = 0$.
Output : Fine-tuned sparse GP with SMD kernel.

- 1 Pretrain $\mathcal{GP}(0, k_{\text{SMD}}(Q_{\text{init}}))$ (initialized by BHI, **Algorithm 2**) M_{init} times;
- 2 Choose $\tilde{\Theta}_{\text{best}}$ with the lowest \mathcal{L} from pretrainings;
- 3 **for** all $k_{\text{SMD}}^{i \times j}$ ($i = j$) **do**
- 4 Obtain the i -th weight w_i from $\tilde{\Theta}_{\text{best}}$;
- 5 **if** $w_i < 1$ **then**
- 6 Remove the i -th component in SMD;
- 7 Remove $\{w_i, \boldsymbol{\mu}_i, \Sigma_i, \boldsymbol{\theta}_i, \boldsymbol{\phi}_i\}$ in $\tilde{\Theta}_{\text{best}}$;
- 8 $Q_{\text{prune}} = Q_{\text{prune}} + 1$;
- 9 **end**
- 10 **end**
- 11 Remove the low intensity dependency structures with $c_{ij} < 1$;
- 12 Fine training the GP with the sparse SMD kernel;

of GPs with the SM kernel. Critically, several de facto inference and learning issues impede the use of the SM kernel, such as hyperparameter initialization and choosing the number of kernel components. The SMD also suffers from these issues. The dependency structures in SMD are dense and therefore need to be sparsified. We propose a structure adaptation (SA) algorithm (see **Algorithm 1**) for the SMD to handle the above issues and to achieve efficient inference and interpretable structure discovery. In **Algorithm 1**, the symbol $\tilde{\Theta}_{\text{best}}$ denotes inferred hyperparameters of better training. Steps 1-2 perform M_{init} trainings to obtain a better hyperparameter position with a smaller loss. Steps 3-10 prune the unimportant components by comparing their weights. Steps 11-12 sparsify the dependency structures by quantifying their intensity, removing the weak ones, and fine-train the GP with sparse dependency structures. In **Algorithm 1**, we use a standard maximum-likelihood approach for the optimization (estimation) of hyperparameters. Such optimization is performed in pretrain stage (step 1) and fine training stage (step 12) of **Algorithm 1**. We have two levels of sparsity for the SMD kernel: the first level of sparsity is obtained from the compression of original components and the second level of sparsity is handled by the reduction of weak dependency structures. Specifically, we introduce the details of the proposed SA algorithm in the following subsections.

5.1. Bootstrap-based hyperparameter initialization (BHI)

The learning of the SMD kernels relies on a good starting point when performing optimization in high-dimensional hyperparameter space. A better initialization can help us more easily discover the underlying structure. Sniffing the structure of the empirical SDs can alleviate the difficulty of hyperparameter initialization due to the connection (indicated by Bochner’s Theorem) between the SD and kernel [8, 17]. However, the empirical SDs are biased estimates

of the true underlying spectral structures, which contain noise and fake peaks denoting spurious patterns. To filter out the noise and fake peaks, we employ the bootstrap techniques [30] to improve the estimation accuracy of the empirical SDs. We draw a large number of spectral samples S^* using bootstrap with replacement from the empirical SDs. We then consider a Gaussian mixture model (GMM) fitting to the bootstrap samples S^* to obtain the Q Gaussians, $p(\tilde{\Theta}|\mathbf{s}) = \sum_{i=1}^Q \tilde{w}_i \mathcal{N}(\mathbf{s}; \tilde{\boldsymbol{\mu}}_i, \tilde{\Sigma}_i)$. Finally, we propose a BHI algorithm (shown in **Algorithm 2**) for the SMD. The bootstrap sampling times $B = 100$ are generally sufficient for robust statistics estimation. The estimates of the hyperparameters obtained in bootstrap are used for initialization. The **Algorithm 2** is performed before optimization in pretrain stage (step 1) of **Algorithm 1**. An illustration of **Algorithm 2** is shown in Fig. 3.

Algorithm 2: Bootstrap-based hyperparameter initialization

Input : $Q_{\text{init}}, B = 100$.

Output : Hyperparameter initialization $\tilde{\Theta}_{\text{init}}$.

- 1 Compute the empirical SD S using the Blackman window and FT;
 - 2 Resample bootstrap spectral samples $S^* \subset S$ from S ;
 - 3 Fit a GMM with Q_{init} components to S^* and obtain a bootstrap estimate $p(\tilde{\Theta}^*|\mathbf{s})$;
 - 4 Sort Q_{init} components with mean position $\tilde{\boldsymbol{\mu}}_i^*$ in $\tilde{\Theta}^*$;
 - 5 Repeat steps 2-4 B times to obtain B estimates $p(\tilde{\Theta}_1^*|\mathbf{s}), p(\tilde{\Theta}_2^*|\mathbf{s}), \dots, p(\tilde{\Theta}_B^*|\mathbf{s})$;
 - 6 The final bootstrap estimates of the hyperparameters are computed as $\tilde{\Theta}_{\text{init}} = \frac{1}{B} \sum_{i=1}^B \tilde{\Theta}_i^*$.
-

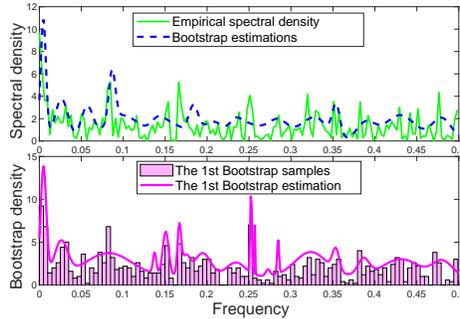


Figure 3: The BHI algorithm on the monthly river flow dataset. The first subplot: the empirical SD (in green line) and bootstrap estimation (in dashed blue line). The second subplot: the 1st bootstrap samples (in magenta bar) and the corresponding estimation (in magenta line). Many small peaks in empirical SD are filtered in the bootstrap estimation.

5.2. Compressed spectral mixture with dependency structures

How to set the number of components is another challenging issue for SM and SMD kernels. We must specify the Q in advance and fix it during optimization. However, inaccurate Q cannot reflect the true number of underlying patterns contained in data, which could lead to overfitting for large Q or underfitting for small Q . This makes all spectral kernels flawed for real-world applications.

To adaptively select the number of components, we first prune the minor components by quantifying their weights. As shown in Fig. 4, we demonstrate the learned importance of components in the SM kernel by using the monthly river flow dataset. Specifically, the weights of components ($i = \{1, 2, 3, 6, 8, 9\}$) in the left subplot are smaller than 1, which means that the corresponding amplitudes in the frequency domain are pretty small. Observing this fact, we simply think a component with a weight smaller than 1 is less important and takes a tiny portion of the signal energy.

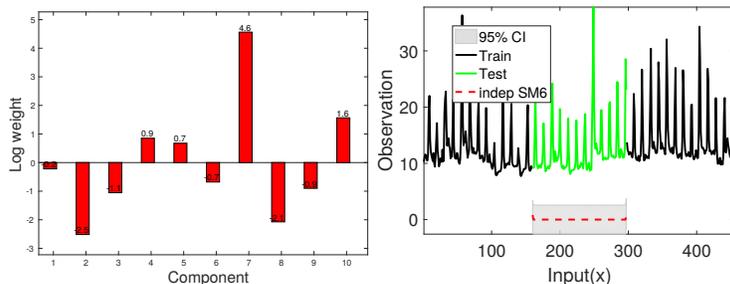


Figure 4: The learned w_i (left) of SM kernel and the interpolation result (right) of the 6-th SM component with a smaller weight. Here, the predictive mean of the 6-th SM component is almost zero and contributes a little to the final predictive distribution and thus can be pruned.

Removing such components with small weights does not affect a GP model’s learning and generalization ability. Consequently, we propose a pruning strategy to compress the SMD. The pruning strategy is described by steps 2-9 of **Algorithm 1**. This strategy can reduce the number of hyperparameters in SMD to $5Q_{\text{rest}}$, where Q_{rest} is the rest components after compression. We define a CR for the SMD, $\alpha_{CR} = 1 - \frac{Q_{\text{rest}}}{Q_{\text{init}}} \times 100\%$, to assess how much the SMD is compressed.

Remark 2. α_{CR} is an indicator of the pruning degree of the SMD using the SA algorithm. In the extreme cases, α_{CR} is %100 if all components are pruned with $Q_{\text{rest}} = 0$ and 100% if all components are kept with $Q_{\text{rest}} = Q_{\text{init}}$.

5.3. Sparse dependency structure and its behavior

In this section, we investigate the sparsity and behavior of the dependency structures in SMD. Observing from Eq. (3) and Eq. (9), there are $Q^2 - Q$ dependency structures. In fact, for two components located far away from each other, the intersection between their SDs is close to zero, which means that their dependency is weak. Eq. (7) indicates that the closer the μ_i , Σ_i and w_i between

components are, the greater the dependency is, and vice versa. Hence, we can confidently remove the low-intensity dependency structures.

Specifically, we introduce a binary mask β_{ij} determined by w_{ij} and a_{ij} to indicate whether remove the dependency structure between components i and j . We have the compressed SMD with sparse dependency structures as

$$k_{\text{SMD}}^{\text{SA}}(\boldsymbol{\tau}) = \sum_{i=1}^{Q_{\text{rest}}} \sum_{j=1}^{Q_{\text{rest}}} \beta_{ij} c_{ij} \exp\left(-\frac{1}{2} \boldsymbol{\tau}_{\theta}^{\top} \Sigma_{ij} \boldsymbol{\tau}_{\theta}\right) \cos(\boldsymbol{\tau}_{\theta}^{\top} \boldsymbol{\mu}_{ij} - \phi_{ij} \pi), \quad (11)$$

where $\beta_{ij} = 0$ if $c_{ij} < 1$ and $i \neq j$; otherwise, $\beta_{ij} = 1$.

Fig. 1 shows the neat contribution of the dependency structure to the final covariance even with zero TP delays. When $\theta \neq 0$ or $\phi \neq 0$, the covariance (in cyan) of the dependency structure is shifted and centered at a different position (shown in subplots (c), (d), and (e) of Fig. 1). We define an SR of the dependency structure for the SMD, $\alpha_{\text{SR}} = \left(1 - \frac{\sum_{i=1}^{Q_{\text{rest}}} \sum_{j=1}^{Q_{\text{rest}}} \beta_{ij}}{Q_{\text{rest}}^2 - Q_{\text{rest}}}\right) \times 100\%$, to evaluate how sparse the dependency structure is.

Remark 3. *The α_{SR} is ensured to be in the range of $[0, 1]$. Note that α_{SR} is 1 if there is no significant dependency structure or 0 if all dependency structures are large.*

6. Experiments

In this section, we comprehensively investigate the performance of the SMD and compare it with that of some state-of-the-art kernels on both synthetic and real-world datasets. For all experiments, the popular kernels implemented in the GPML toolbox [3] are used as baselines, such as the linear (LIN), SE, polynomial (Poly), periodic (PER), rational quadratic (RQ), Matérn (MA), Gabor, fractional Brownian motion covariance (FBM), underdamped linear Langevin process covariance (ULL), neural network (NN) and SM kernels. The same number of components Q is used for the SM, NSM, and SMD kernels. In all plots, the training data, testing data, SM prediction, SMD prediction, and CI are shown in black, green, red, blue, and gray shadow, respectively.

6.1. Model assessment

We consider multiple metrics to assess the performances and characteristics of GP models, such as

- the mean squared error (MSE) defined as $\text{MSE} \triangleq \frac{1}{n} \sum_{i=1}^n (y_i^* - \tilde{y}_i^*)^2$ to measure the generalization performance of GP;
- the 95% CI (instead of, e.g., error bar) to visualize the uncertainty of prediction;
- the posterior correlation ρ_{ij} to quantify the latent dependency between SM components;

- the γ_{ij} (see Eq. (10)) to indicate the intensity of dependency structure learned by the SMD;
- the CR (α_{CR}) and SR (α_{SR}) of the SMD.

6.2. Learning a synthetic signal with dependency structure

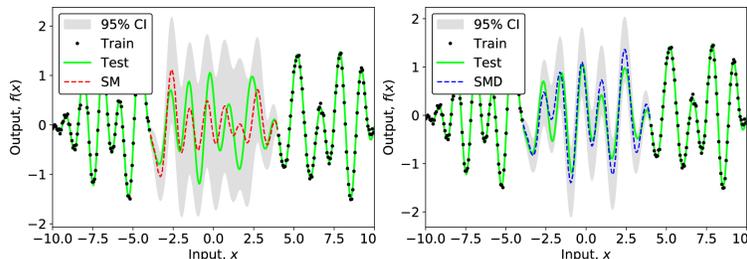


Figure 5: Performance of the SM (left) and SMD (right) on a synthetic signal.

We illustrate the capability of the SMD to capture TP delayed dependency structure in a synthetic signal. The signal is sampled from the following GP with a hybrid kernel structure: $f(x) \sim \mathcal{GP}(0, k_{\text{SMD}}(\theta = \{0.1, 0.3\}, \phi = \{0.1, 0.3\}) + k_{\text{SM}})$. The signal contains a dependency structure due to the employment of k_{SMD} . The SM and SMD kernels of the signal $f(x)$ have $Q = 2$ components. We generate a time series of length 300 in the interval $[-10, 10]$ and add some noise to it (see Fig. 5). In this experiment, we remove the middle 40% of the signal and consider it as missing testing data (in green). The rest of the signal forms the training data (in black). Both SMD and SM are configured with $Q = 5$ and with the same initial values of the hyperparameters w_i, μ_i, σ_i^2 . Other hyperparameters of SMD, θ_i and ϕ_i , are initialized to be zeros.

In Fig. 5, the performance difference between SMD and SM is clear in terms of the mean and uncertainty of the prediction. The SMD is capable of learning the hybrid covariance with dependency structure well. For the SM (dashed red), it is more difficult to recognize such a dependency structure and to interpolate the missing block. Here, the SMDs without TP delays, with only time delay, and with only phase delay cannot interpolate the missing block well. Obviously, the SMD yields better prediction and smoother CIs than the SM and therefore achieves the lowest MSE (see Table 2).

6.3. Long range interpolation of monthly river flow monitoring

Interpolation is a well-known task for GP learning. In this experiment, we validate the long-range interpolation ability of GP with the SMD kernel. We consider the monthly river flow dataset because it reveals time and phase patterns with variability. The moon and sun are primarily responsible for the rising and falling of tidal river flows, and such effects are delayed and augmented by gravity and resonances. The mean monthly river flow in the Madison River near West Yellowstone is the average flow from 1923 to 1960 [31]. Empirical analysis [31]

shows various characteristics of this flow data experiment: short term monthly variations, medium term seasonal patterns, irregular periodic long term trends caused by the relative positions of the moon and sun, and some white noise.

As such, the monthly river flow contains complicated patterns (see Fig. 6) that may be caused by physical interferences. The appearance time of the flow peak is not periodical, and its amplitude is always irregular. There are 456 records in the dataset. Here, 30% of the data, namely, the long range middle part, is removed for testing, while the rest of the data are used for training.

In Table 2 and Fig. 6, the results indicate that both the SMD and SM can interpolate the missing month river flow well. However, the SMD achieves better performance and confidence. The SMD is generally more effective in modeling complex patterns hidden in these data. Furthermore, using the same initial Q and the SA algorithm, Table 3 shows that SMD achieves a better CR of 38.9% than SM. The SR of 89.3% indicates that most of low intensity dependency structures in the SMD are removed.

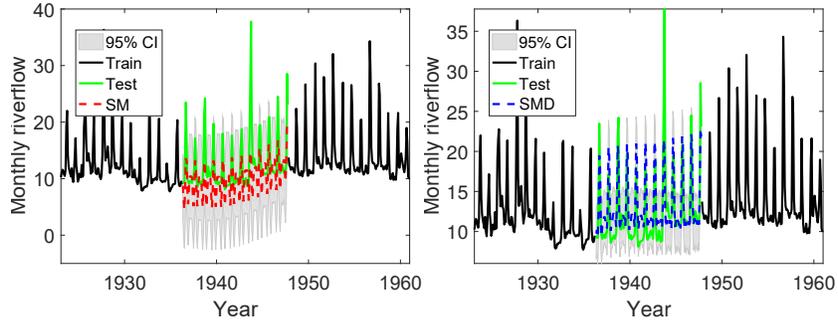


Figure 6: Interpolations of GPs with the SM (left) and SMD (right) kernels on the monthly river flow dataset.

Posterior dependency between SM components: On the other hand, by computing the posterior covariance between two functions, conditioned on their sum [21], we obtain posterior cross covariance between two components as, $\text{Cov}(f_{SM,i}^*, f_{SM,j}^* | \mathbf{f}_{SM,i+j}) = -\mathbf{k}_{SM,i}^{*\top} K_{SM,i+j}^{-1} \mathbf{k}_{SM,j}^*$, where $\mathbf{f}_{SM,i+j} = \mathbf{f}_{SM,i} + \mathbf{f}_{SM,j}$ and $K_{SM,i+j} = K_{SM,i} + K_{SM,j}$. As investigated in [21], the posterior cross covariance can indicate the underlying dependency between $f_{SM,i}$ and $f_{SM,j}$. We further normalize the posterior cross covariance as posterior correlation coefficient ρ_{ij} with range $[-1, 1]$: $\rho_{ij} = \frac{\text{Cov}(f_{SM,i}^*, f_{SM,j}^* | \mathbf{f}_{SM,i+j})}{(\mathbb{V}(f_{SM,i}^* | \mathbf{f}_{SM,i+j}) \mathbb{V}(f_{SM,j}^* | \mathbf{f}_{SM,i+j}))^{1/2}}$. Note that there is no underlying dependency if $\rho_{ij} = 0$, otherwise, $f_{SM,i}$ and $f_{SM,j}$ are dependent. In Fig. 7, the left subplot shows high and complex posterior correlation coefficient ρ_{56} of the SM. In the right subplot, the intensity of the dependency structure in the SMD is indicated by the positive and negative values of γ_{14} . Note that there is no alignment between SM and SMD components because they are separately optimized.

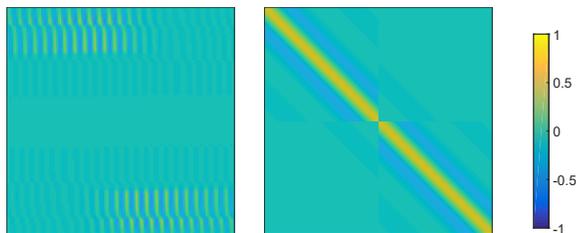


Figure 7: The posterior correlations ρ_{ij} (left) and the dependency structure intensity γ_{ij} (right) for the SM and SMD kernels on river flow, respectively.

6.4. Dependency structure for joint interpolation and extrapolation of yearly sunspots modeling

In addition to the interpolation task, we simultaneously perform interpolation and extrapolation to further substantiate the learning ability of the SMD. We consider the yearly sunspot number dataset¹ collected between 1700 and 2014. The historical evolution of yearly sunspots can help explain spatial magnetic field environment changes affected by sun activities. The yearly sunspot number is obtained by taking an arithmetic mean of the daily total sunspot number over all days of each year. Sunspots appear darker than the surrounding areas on the sun’s photosphere [32]. They usually have lower surface temperatures than the areas around them. A sunspot has an irregular period of existence: the average number of sunspots that are monitored increases and decreases with a quasi-period. There are dependencies between patterns of sunspots caused by some physical types of interference. As shown in Fig. 8, patterns in yearly sunspots contain various irregular peaks over 315 years.

There are 315 records in the dataset. We use the last 10% of the data for the extrapolation test (in solid green) and randomly sample 20% of the original data from the first 90% of the yearly sunspots as the interpolation test (in crossed green). The remaining 70% of the data are used for training (in black). The legends Test_{ext} and Test_{int} denote the extrapolation and interpolation testing data, respectively. Note that the training data are not equally sampled due to missing values. We initially considered $Q = 20$ components for both SMD and SM. Time and phase delays in the SMD are also initialized as zeros.

As shown in subplots (a) and (b) of Fig. 9, there is a significant TP delay dependency structure (in cyan) between components 2 (in green) and 4 (in black) in the SMD. As investigated in Section 5.3, $\hat{k}_{\text{SMD}}^{2 \times 4}$ is significant because component 4 and component 2 are close to each other. Due to the TP delays, the covariance of $k_{\text{SMD}}^{2 \times 4}$ is shifted to left. The period of the dependency structure $k_{\text{SMD}}^{2 \times 4}$ is smaller than the 2nd component and larger than the 4th component and has a connection to the location of $\hat{k}_{\text{SMD}}^{2 \times 4}$. From the CR and SR shown in Table 3, the SMD using the SA algorithm can reduce the hyperparameter

¹<http://www.sidc.be/silso/infosntotyearly>

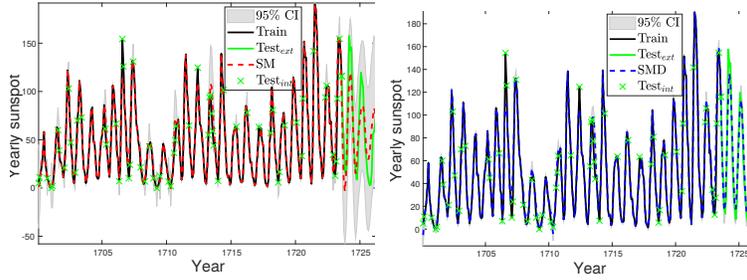


Figure 8: Interpolations and extrapolations of GPs with the SM (left) and SMD (right) kernels on the yearly sunspot dataset.

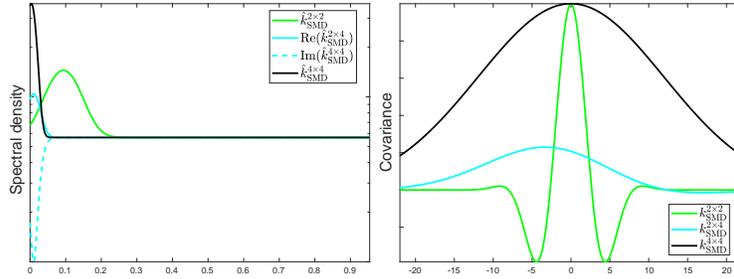


Figure 9: The learned dependency structures on the yearly sunspot dataset. Left subplot: the SDs of $\hat{k}_{\text{SMD}}^{2 \times 2}$, real part $\text{Re}(\hat{k}_{\text{SMD}}^{2 \times 4})$, imaginary part $\text{Im}(\hat{k}_{\text{SMD}}^{2 \times 4})$, and $\hat{k}_{\text{SMD}}^{4 \times 4}$ are in green, solid cyan, dashed cyan, and black, respectively. Right subplot: the covariances of $k_{\text{SMD}}^{2 \times 2}$, $k_{\text{SMD}}^{2 \times 4}$, and $k_{\text{SMD}}^{4 \times 4}$.

size by 40.7% and dependency structures by 50.3%. The above results indicate that both the SMD and SM can interpolate missing values well with a small CI. However, for the extrapolation task, the SMD achieves better performance and confidence (see Table 2 and Fig. 8). With this experimental result, we can conclude that the SMD can perform interpolation and extrapolation equally well for incomplete signals.

6.5. Scalable SMD on large scale multidimensional data

Furthermore, we comparatively evaluate the scalable SMD on a large multidimensional abalone² dataset. We apply automatic relevance determination (ARD) for other baseline kernels to remove irrelevant input. Note that the FBM, ULL, and NSM kernels are not applicable to multidimensional datasets ($P > 2$). When modeling large data, exact inference [15, 33, 34, 35] of GP is prohibitively expensive and meets $\mathcal{O}(n^3)$ computational complexity and $\mathcal{O}(n^2)$ memory. The expensive computation cost involves computing the inverse and determinant of $K + \sigma_n^2 I$. We consider a scalable SMD using stochastic variational inference

²<http://archive.ics.uci.edu/ml/datasets/abalone>

(SVI) framework [36, 37]. Specifically, SVI can approximate the underlying GP posterior with a GP conditioned on a small set U with m inducing points. The inducing points U can be seen as a set of global variables summarizing the structure of the large training data to perform variational inference. The variational distribution $p(\mathbf{u}) = \mathcal{N}(\mathbf{u}; \boldsymbol{\mu}_{\mathbf{u}}, \Sigma_{\mathbf{u}})$ with

$$\begin{aligned}\Sigma_{\mathbf{u}} &= K_{UU}^{-1} + \sigma_n^{-2} K_{UU}^{-1} K_{UX} K_{UX}^{\top} K_{UU}^{-1}, \\ \boldsymbol{\mu}_{\mathbf{u}} &= \sigma_n^{-2} \Sigma_{\mathbf{u}}^{-1} K_{UU}^{-1} K_{UX} \mathbf{y},\end{aligned}\tag{12}$$

gives a variational lower bound $\mathcal{L}_3(\mathbf{u}; \boldsymbol{\mu}_{\mathbf{u}}, \Sigma_{\mathbf{u}})$, also called evidence lower bound (ELBO) of the quantity $p(\mathbf{y}|X)$, satisfying $\log p(\mathbf{y}|X) \geq \mathcal{L}_3(\mathbf{u}; \boldsymbol{\mu}_{\mathbf{u}}, \Sigma_{\mathbf{u}})$, where $K_{UX} = k_{\text{SMD}}(U, X)$ and $K_{UU} = k_{\text{SMD}}(U, U)$. As mentioned in [36], the variational distribution $p(\mathbf{u})$ contains all the information encoded into the posterior approximation, which describes the distribution of function values at the inducing points U . Letting $\frac{\partial \mathcal{L}_3}{\partial \boldsymbol{\mu}_{\mathbf{u}}} = 0$ and $\frac{\partial \mathcal{L}_3}{\partial \Sigma_{\mathbf{u}}} = 0$, we can approximate the optimal solution of the variational distribution. Therefore, we have the posterior distribution of a test point as $p(y^*|\mathbf{x}^*, X, Y) = \mathcal{N}(\tilde{y}^*, \mathbb{V}[y^*])$, where the predictive mean $\tilde{y}^* = \mathbf{k}_U^* K_{UU}^{-1} \boldsymbol{\mu}_{\mathbf{u}}$, predictive variance $\mathbb{V}[y^*] = k^{**} + \mathbf{k}_U^{*\top} (K_{UU}^{-1} \Sigma_{\mathbf{u}} K_{UU}^{-1} - K_{UU}^{-1}) \mathbf{k}_U^*$, $k^{**} = k_{\text{SMD}}(\mathbf{x}^*, \mathbf{x}^*)$, and $\mathbf{k}_U^* = k_{\text{SMD}}(U, \mathbf{x}^*)$. Finally, the complexity of the SMD on a large dataset is reduced to $\mathcal{O}(m^3)$.

The abalone³ dataset has 4177 samples with eight attributes: sex, length, diameter, height, whole weight, shucked weight, visceral weight, and shell weight. We aim to predict the age of abalone, which is usually measured by physical assessment. Generally, an abalone’s age is measured by cutting the shell through the cone, staining it, and counting the number of rings through a microscope. The number of rings directly reflects the age of an abalone. The task is to predict the number of rings from the eight attributes. Specifically, we use the first 3377 instances as training data and the remaining 800 instances as testing data. We use $Q = 10$ components for the SMD and SM due to the explosive expansion of hyperparameter space for multidimensional input. Note that we set the number of inducing points $m = 500$ in the SVI. The SMD and SM use the same hyperparameters initialization described in the SA algorithm. The results in Table 2 show that on this type of task, the SMD also performs better, with a lower MSE than the SM.

6.6. Discussion

In Table 2, the results demonstrate that the SMD ($\boldsymbol{\theta} \neq \mathbf{0}, \boldsymbol{\phi} \neq \mathbf{0}$) performs better than other baselines as well as the other SMD variants ($\boldsymbol{\theta} = \mathbf{0}$ or $\boldsymbol{\phi} = \mathbf{0}$). Our experiment analysis indicates that signals containing dependency structure caused by physical interference can be learned well by a GP with SMD kernel. The SM and SMDs using BHI perform much better than those using random initialization (RI). The proposed BHI in the SA algorithm can provide a good starting point for optimization in high-dimensional hyperparameter space.

³<http://archive.ics.uci.edu/ml/datasets/abalone>

Table 2: Performances of various baselines versus the SMD in terms of MSE. RI denotes random initialization (RI).

Kernel	Synthetic	Riverflow	Sunspot	Abalone
LIN	0.32±0.11	24.37±4.57	1617.45±131.63	10.93±2.07
SE	0.31±0.20	174.14±19.23	591.24±9.21	8.14±3.24
Poly	0.32±0.18	182.01±18.73	1627.02±127.84	6.30±1.71
PER	0.35±0.13	20.31±7.22	1533.84±143.15	7.98±2.80
RQ	0.31±0.10	24.88±6.32	307.22±30.38	5.38±1.53
MA	0.32±0.14	170.29±18.90	590.72±26.21	7.52±1.05
Gabor	0.31±0.20	22.51±3.37	4047.83±19.78	3.68±1.13
FBM	0.49±0.19	23.84±3.21	6428.69±515.98	--
ULL	0.26±0.11	168.06±18.16	467.08±23.90	--
NN	0.32±0.16	23.65±2.32	1522.67±68.58	3.61±1.07
NSM	0.53±0.21	166.85±20.45	4697.24±375.15	--
SM (RI)	0.41±0.16	23.55±2.79	363.64±64.32	3.64±1.12
SM (BHI)	0.43±0.18	13.91±1.78	270.78±13.96	3.52±0.51
SMD (RI, $\theta = \mathbf{0}, \phi = \mathbf{0}$)	0.47±0.13	23.31±4.76	329.81±18.19	3.56±1.90
SMD (BHI, $\theta = \mathbf{0}, \phi = \mathbf{0}$)	0.28±0.14	9.55±1.79	184.47±18.58	3.37±0.42
SMD (RI, $\theta \neq \mathbf{0}, \phi \neq \mathbf{0}$)	0.45±0.95	19.26±3.53	322.50±21.67	3.74±1.31
SMD (BHI, $\theta \neq \mathbf{0}, \phi \neq \mathbf{0}$)	0.05±0.01	8.96±0.72	171.16±13.10	3.25±0.38

Table 3: The average CRs and SRs of the SM and SMD variants using SA.

Kernel	CR Riverflow	CR Sunspot	SR Riverflow	SR Sunspot
SM	33.2%	30.3%	–	–
SMD (BHI, $\theta = \mathbf{0}, \phi = \mathbf{0}$)	35.3%	32.8%	89.1%	57.6%
SMD (BHI, $\theta \neq \mathbf{0}, \phi \neq \mathbf{0}$)	38.9%	40.7%	89.3%	55.3%

As shown in Table 3, the CRs of both the SM and SMD are larger than 30% due to the use of the SA algorithm. Hence, the SA algorithm can much reduce the hyperparameter space of SMD and SM. In Table 3, the SMDs with SA usually have better CR than the SM. This may be caused by the fact that the SMD with dependency structures has a better representation ability than the SM. The SMD can describe an underlying function with fewer components than the SM because the latter needs additional components to delineate the latent dependency structure. In addition, Table 3 shows that the SR of all the SMDs is high, which means that the dependency structures are sparse. Most dependency structures are tiny and removable. Due to the higher number of hyperparameters (at least three times than SM) and overfitting troubles, the results in Table 2 show the unsatisfactory performance of NSM on the synthetic signal and on real-world datasets.

7. Conclusion

We propose a novel SMD kernel, which extends the SM kernel by incorporating TP delayed dependency structures. An interpretable SA algorithm for the SMD

is introduced to effectively initialize its hyperparameters, compress components, and obtain sparse dependency structures automatically.

The results of extensive experiments on both the synthetic and real-life datasets indicate that the SMD using the structure adaptation (SA) algorithm can learn TP delayed dependency structures between the components and perform more accurate interpolation and extrapolation. Hence, the benefits of SMD are shown to be significant.

Two main issues remain to be addressed in future work. The first issue is the initialization of the TP parameters. Here, we simply initialized them as zeros. However, more tailored, effective methods remain to be investigated. Another issue, common to all GP methods, is the problem of sparse or efficient inference [3, 34], which also needs to be further improved for GPs with the SMD on big data.

Acknowledgment

We would like to thank Elena Marchiori, Twan van Laarhoven, and Perry Groot for their comments on a past version of this work. This research is supported by the Guangdong Provincial Key Laboratory of Future Networks of Intelligence, The Chinese University of Hong Kong, Shenzhen, under Grant No. 2022B1212010001. The work was partly supported by the Natural Science Foundation of China (NSFC) with grant No. 62106212, by the Natural Science Foundation of Hunan Province, China, under Grant 2023JJ40689, and by the High Performance Computing Center of Central South University (CSU). The work of Feng Yin was supported by the NSFC with Grant No. 62271433.

References

- [1] S. Theodoridis, *Machine learning: a Bayesian and optimization perspective*, 2nd edition, Academic press, 2020.
- [2] C. E. Rasmussen, H. Nickisch, Gaussian processes for machine learning (GPML) toolbox, *Journal of Machine Learning Research* 11 (Nov) (2010) 3011–3015.
- [3] C. E. Rasmussen, C. K. Williams, *Gaussian process for machine learning*, MIT press, 2006.
- [4] F. Yin, F. Gunnarsson, Distributed recursive Gaussian processes for rss map applied to target tracking, *IEEE Journal of Selected Topics in Signal Processing* 11 (3) (2017) 492–503.
- [5] Y. Xu, F. Yin, W. Xu, J. Lin, S. Cui, Wireless traffic prediction with scalable Gaussian process: Framework, algorithms, and verification, *IEEE J. Sel. Areas Commun.* 37 (6) (2019) 1291–1306.

- [6] F. Yin, Z. Lin, Q. Kong, Y. Xu, D. Li, S. Theodoridis, S. R. Cui, Fedloc: Federated learning framework for data-driven cooperative localization and location data processing, *IEEE Open Journal of Signal Processing* 1 (2020) 187–215.
- [7] K. Chen, Q. Kong, Y. Dai, Y. Xu, F. Yin, L. Xu, S. Cui, Recent advances in data-driven wireless communication using Gaussian processes: A comprehensive survey, *China Communications* 19 (2022) 218–237.
- [8] A. Wilson, R. Adams, Gaussian process kernels for pattern discovery and extrapolation, in: *Proceedings of the 30th International Conference on Machine Learning (ICML-13)*, 2013, pp. 1067–1075.
- [9] D. Duvenaud, J. R. Lloyd, R. Grosse, J. B. Tenenbaum, Z. Ghahramani, Structure discovery in nonparametric regression through compositional kernel search, *arXiv preprint arXiv:1302.4922*.
- [10] F. Yin, L. Pan, T. Chen, S. Theodoridis, Z.-Q. T. Luo, A. M. Zoubir, Linear multiple low-rank kernel based stationary Gaussian processes regression for time series, *IEEE Transactions on Signal Processing* 68 (2020) 5260–5275.
- [11] Y. Dai, T. Zhang, Z. Lin, F. Yin, S. Theodoridis, S. Cui, An interpretable and sample efficient deep kernel for Gaussian process, in: *Conference on Uncertainty in Artificial Intelligence*, PMLR, 2020, pp. 759–768.
- [12] R. Dürichen, M. A. Pimentel, L. Clifton, A. Schweikard, D. A. Clifton, Multitask Gaussian processes for multivariate physiological time-series analysis, *IEEE Transactions on Biomedical Engineering* 62 (1) (2015) 314–322.
- [13] M. Kupilik, F. Witmer, E.-A. MacLeod, C. Wang, T. Ravens, Gaussian process regression for arctic coastal erosion forecasting, *arXiv preprint arXiv:1712.00867*.
- [14] K. Chen, T. van Laarhoven, P. Groot, J. Chen, E. Marchiori, Multioutput convolution spectral mixture for Gaussian processes, *IEEE Transactions on Neural Networks and Learning Systems*.
- [15] A. G. Wilson, E. Gilboa, A. Nehorai, J. P. Cunningham, Fast kernel learning for multidimensional pattern extrapolation, in: *Advances in Neural Information Processing Systems*, 2014, pp. 3626–3634.
- [16] S. Remes, M. Heinonen, S. Kaski, Non-stationary spectral kernels, in: *Advances in Neural Information Processing Systems*, 2017, pp. 4645–4654.
- [17] W. Herlands, A. Wilson, H. Nickisch, S. Flaxman, D. Neill, W. Van Panhuis, E. Xing, Scalable Gaussian processes for characterizing multidimensional change surfaces, in: *Artificial Intelligence and Statistics*, 2016, pp. 1013–1021.

- [18] K. Chen, T. van Laarhoven, J. Chen, E. Marchiori, Incorporating dependencies in spectral kernels for Gaussian processes, in: Machine Learning and Knowledge Discovery in Databases - European Conference, ECML PKDD 2019, Würzburg, Germany, 2019, Proceedings, 2019.
- [19] S. Bochner, Lectures on Fourier Integrals.(AM-42), Vol. 42, Princeton University Press, 2016.
- [20] M. L. Stein, Interpolation of spatial data: some theory for Kriging, Springer Science & Business Media, 2012.
- [21] D. Duvenaud, Automatic model construction with Gaussian processes, Ph.D. thesis, University of Cambridge (2014).
- [22] A. G. Wilson, Covariance kernels for fast automatic pattern discovery and extrapolation with Gaussian processes, University of Cambridge.
- [23] K. Chen, T. van Laarhoven, E. Marchiori, Gaussian processes with skewed Laplace spectral mixture kernels for long-term forecasting, Machine Learning 110 (8) (2021) 2213–2238.
- [24] P. A. Jang, A. Loeb, M. Davidow, A. G. Wilson, Scalable Lévy process priors for spectral kernel learning, in: Advances in Neural Information Processing Systems, 2017, pp. 3943–3952.
- [25] A. Klenke, Wahrscheinlichkeitstheorie, Springer Spektrum Berlin, Heidelberg, 2006.
- [26] M. Loeve, Probability theory, Courier Dover Publications, 2017.
- [27] G. Gaspari, S. E. Cohn, Construction of correlation functions in two and three dimensions, Quarterly Journal of the Royal Meteorological Society 125 (554) (1999) 723–757.
- [28] M. G. Genton, W. Kleiber, et al., Cross-covariance functions for multivariate geostatistics, Statistical Science 30 (2) (2015) 147–163.
- [29] H. Bateman, Tables of integral transforms [volumes I & II], Vol. 1, McGraw-Hill Book Company, 1954.
- [30] A. M. Zoubir, B. Boashash, The Bootstrap and its application in signal processing, IEEE signal processing magazine 15 (1) (1998) 56–76.
- [31] K. W. Hipel, A. I. McLeod, Time series modelling of water resources and environmental systems, Vol. 45, Elsevier, 1994.
- [32] R. C. Willson, S. Gulkis, M. Janssen, H. S. Hudson, G. A. Chapman, Observations of solar irradiance variability, Science 211 (4483) (1981) 700–702. [doi:10.1126/science.211.4483.700](https://doi.org/10.1126/science.211.4483.700).

- [33] C. K. Williams, M. Seeger, Using the Nyström method to speed up kernel machines, in: *Advances in neural information processing systems*, 2001, pp. 682–688.
- [34] J. Quiñero-Candela, C. E. Rasmussen, A unifying view of sparse approximate Gaussian process regression, *Journal of Machine Learning Research* 6 (Dec) (2005) 1939–1959.
- [35] E. Snelson, Z. Ghahramani, Sparse Gaussian processes using pseudo-inputs, in: *Advances in neural information processing systems*, 2006, pp. 1257–1264.
- [36] J. Hensman, N. Fusi, N. D. Lawrence, Gaussian processes for big data, in: *Proceedings of the Twenty-Ninth Conference on Uncertainty in Artificial Intelligence, UAI 2013, Bellevue, WA, USA, August 11-15, 2013*, pp. 282–290.
- [37] J. Hensman, N. Durrande, A. Solin, Variational Fourier features for Gaussian processes, *Journal of Machine Learning Research* 18 (2017) 151:1–151:52.