

Linear Multiple Low-Rank Kernel Based Stationary Gaussian Processes Regression for Time Series

Feng Yin, Lishuo Pan, Xinwei He, Tianshi Chen, Sergios Theodoridis, Zhi-Quan (Tom) Luo
 School of Science and Engineering, The Chinese University of Hong Kong, Shenzhen
 Shenzhen Research Institute of Big Data (SRIBD)
 Longxiang blvd 2001, Longgang district, Shenzhen, China, 518172.

arXiv:1904.09559v1 [cs.LG] 21 Apr 2019

Abstract—Gaussian processes (GP) for machine learning have been studied systematically over the past two decades and they are by now widely used in a number of diverse applications. However, GP kernel design and the associated hyper-parameter optimization are still hard and to a large extent open problems. In this paper, we consider the task of GP regression for time series modeling and analysis. The underlying stationary kernel can be approximated arbitrarily close by a new proposed grid spectral mixture (GSM) kernel, which turns out to be a linear combination of low-rank sub-kernels. In the case where a large number of the sub-kernels are used, either the Nyström or the random Fourier feature approximations can be adopted to deal efficiently with the computational demands. The unknown GP hyper-parameters consist of the non-negative weights of all sub-kernels as well as the noise variance; their estimation is performed via the maximum-likelihood (ML) estimation framework. Two efficient numerical optimization methods for solving the unknown hyper-parameters are derived, including a sequential majorization-minimization (MM) method and a non-linearly constrained alternating direction of multiplier method (ADMM). The MM matches perfectly with the proven low-rank property of the proposed GSM sub-kernels and turns out to be a part of efficiency, stable, and efficient solver, while the ADMM has the potential to generate better local minimum in terms of the test MSE. Experimental results, based on various classic time series data sets, corroborate that the proposed GSM kernel-based GP regression model outperforms several salient competitors of similar kind in terms of prediction mean-squared-error and numerical stability.

Index Terms—ADMM, Gaussian processes, hyper-parameter optimization, majorization-minimization, linear multiple kernel, low-rank kernels, prediction, time series.

I. INTRODUCTION

Gaussian processes (GP) constitute a class of important Bayesian non-parametric models for machine learning and they are tightly connected to several other popular models, such as support vector machines (SVM), regularized-least-squares, relevance vector machines and auto-regressive-moving-average (ARMA), single-layer Bayesian neural networks [2] and, more recently, to the deep neural networks [3], [4]. Gaussian processes are also used as outstanding surrogate functions for Bayesian optimization nowadays [5]. The idea behind the GP models is to impose a Gaussian prior on the underlying function/system and then compute the posterior distribution over the function given the observed data. GP

models have been used in a plethora of applications due to their outstanding performance in function approximation with a natural uncertainty bound.

Gaussian processes models are simple in terms of mathematical formulation and analysis thanks to the underlying Gaussian assumption. However, like other kernel methods, such as support vector machines, one major problem with GP models lies in the selection of an appropriate kernel function. It is well known that a good kernel function is capable of lifting the raw features to a much higher (even infinite) dimensional space, where regression and classification can be done more effectively, e.g., [6]. In practice, kernel selection is often done subjectively, relying on eye inspection of data patterns and a handful of elementary kernels such as the linear kernel, squared-exponential (SE) kernel, Matérn kernel and their hybrid are popular alternatives. For instance, the SE kernel was used for sport trajectory modeling in [7] and for financial data modeling and prediction in [8], while linear and Matérn kernels were used for energy load forecasting in [9], to mention a few in different sectors, even though the selected kernel may not fit the data well.

In order to bypass the need for human intervention, automatic and optimal kernel design is largely demanded. One option is to resort to multiple kernel learning techniques. Multiple kernel refers to learning a linear or nonlinear combination of primitive kernels systematically for a target machine learning (supervised, un-supervised, etc.) model, via a specific optimization method with the goal to let data determine the best kernel configuration. This idea has been implemented mostly based on linear multiple kernel (LMK) for supervised SVM models [10], [11], for supervised regularized least-squares models [12], and for un-supervised data clustering [13], etc. The idea of mixing elementary kernels for Gaussian process regression also exists, e.g., for CO₂ prediction in [2] and for other prediction tasks in a few recent works [9], [14], [15]. However, the main drawback is that the primitive kernels are selected subjectively and primitively combined with equal weights. In other words, the weights were pre-selected and not learnt via optimizing a performance metric, and the resulting simple equal-weighted linear combination of primitive kernels may be way sub-optimal for fitting the given data.

There also exist some competing universal kernel design methods. In [16], Lazaro-Gredilla *et.al.* proposed a sparse spectrum Gaussian process (SSGP) that extends the linear trigonometric Bayesian model. The spectral density of a

The conference version of this paper [1] has been published in the proceedings of 21st International Conference on Information Fusion (FUSION), University of Cambridge, Cambridge, UK, July, 2018.

stationary covariance kernel is sparsified to approximate the standard GP. The SSGP learns the model hyper-parameters, including the spectral points, precision of a prior, noise variance as well as the lengthscales of the automatic relevance determination (ARD) kernel via maximizing the marginal likelihood with the conjugate gradient method. In [17], Duvenaud *et.al.* defined a space of kernel structures built compositionally by adding and multiplying a small number of primitive kernels and search for the optimal combination over the space. In [18], [19], Wilson *et.al.* proposed a spectral mixture (SM) kernel with the idea to approximate the spectral density with a Gaussian mixture model first in the frequency domain and transform it back into the time domain.

The predictive performance of GP regression depends on the goodness of the model parameters, often referred to as hyper-parameters. There exist two classes of methods for tuning the GP hyper-parameters. The class of deterministic methods consists of the maximum likelihood (ML) estimation based method and cross-validation (CV) based method among others [2], [20]. The class of stochastic methods includes for instance the hybrid Monte-Carlo and Markov chain Monte-Carlo (MCMC) sampling methods [21], [22]. In this paper, we follow the deterministic ML based method that is more widely used in the GP community.

The main contributions of this work are the following:

- Based on the assumption that there exists a true kernel and moreover it is stationary, we propose a novel grid spectral mixture (GSM) kernel for time series modeling and analysis. The GSM kernel simply modifies the original spectral mixture kernel [19] by fixing the frequency and variance parameters to a set of pre-selected grids while leaving only the weights to be optimized.
- As a major contribution, the resulting GSM kernel belongs to the class of linear multiple kernels, and the associated sub-kernels are proven to have low-rank property under reasonable conditions. Moreover, by fixing the grids, the ML based hyper-parameter optimization task becomes equivalent to a difference-of-convex problem with nicer structure to be dealt with. When the proposed GSM kernel contains a large number of sub-kernels, we propose to apply Nyström or random Fourier feature approximation for saving in computational complexity and storage requirements.
- As another major contribution, we derive two effective numerical methods for tuning the GP hyper-parameters. The first method is a sequential majorization-minimization (MM) method, and the second one is a non-linearly constrained alternating direction of multiplier method (ADMM). The former method turns out to be very fast and stable, while the latter method has the potential to achieve a better local minimum in the sense of achieving smaller prediction mean-squared-error (MSE). For both methods, the solution turns out to be sparse, which is a welcome feature in the context of data overfitting problem.
- Tests based on eight standard time series data sets in various aspects verify that the proposed GSM kernel for GP modeling, empowered with an efficient hyper-parameter

optimization approach, is able to achieve much improved prediction performance and robustness as compared to other competing GP models of similar kind.

The remainder of this paper is organized as follows. Section II provides the background about Gaussian process regression, the classic ML based hyper-parameter optimization, and the linear multiple kernel. Section III first reviews the SM kernel, followed by a new GSM kernel, which turns out to be a linear multiple kernel. Section IV introduces the Nyström and random Fourier feature approximation of the GSM sub-kernels for computational and memory savings. Section V first presents the ML based hyper-parameter estimation problem for large scale linear multiple kernel, including the proposed GSM kernel and it further presents two numerical optimization methods, namely a sequential MM method and an ADMM method. Experimental results are given in Section VI. Finally, Section VII concludes this paper. Proofs of some important properties of the GSM kernel are given in Appendix.

Notation: Throughout this paper, matrices are presented with boldface uppercase letters, vectors with boldface lowercase letters, and scalars with normal lowercase letters. We use \mathbb{R} to denote the set of real numbers. The operator $[\cdot]^T$ stands for vector/matrix transpose, $\text{tr}(\cdot)$ for trace of a square matrix, $\text{rank}(\cdot)$ for rank of a matrix, $\|\cdot\|_p$ for L_p norm of a vector and $\|\cdot\|_F$ for the Frobenius norm of a matrix, $\mathbb{E}_{p(x)}(\cdot)$ for the expectation taken with respect to the probability density function (PDF) $p(x)$, ∇_{θ} for gradient, $\mathcal{N}(v; \mu, \sigma^2)$ for Gaussian distribution of a random variable V with mean μ and variance σ^2 , $\det(\cdot)$ is determinant of a matrix, $\text{erf}(x)$ is Gaussian error function, $[a]_+$ takes the maximum between a and zero. Lastly, $\mathbf{X} \succeq \mathbf{Y}$ means $\mathbf{X} - \mathbf{Y}$ is positive semi-definite, $\langle \mathbf{X}, \mathbf{Y} \rangle$ is the inner product of two square matrices, $\mathbf{X} \circ \mathbf{Y}$ stands for the Hadamard (point-wise) matrix multiplication of \mathbf{X} and \mathbf{Y} .

II. BACKGROUND

In this section, we first review GP regression in subsection II-A and classic ML based GP hyper-parameter optimization in subsection II-B. Lastly, we introduce linear multiple kernel in subsection II-C.

A. GP Regression

A Gaussian process is a collection of random variables, any finite subset of which follows a Gaussian distribution [2]. In the sequel, we solely focus on scalar output, real-valued Gaussian processes that are completely specified by a mean function and a kernel function (a.k.a. covariance function). Concretely, we express

$$f(\mathbf{x}) \sim \mathcal{GP}(m(\mathbf{x}), k(\mathbf{x}, \mathbf{x}'; \boldsymbol{\theta}_h)), \quad (1)$$

where $m(\mathbf{x})$ is the mean function, which is often set to zero in practice, especially when there is no prior knowledge available; and $k(\mathbf{x}, \mathbf{x}'; \boldsymbol{\theta}_h)$ is the kernel function tuned by the kernel hyper-parameters, $\boldsymbol{\theta}_h$.

Let us consider the following GP regression model

$$y = f(\mathbf{x}) + e, \quad (2)$$

where $y \in \mathbb{R}$ is a continuous-valued, scalar output; the unknown function $f(\mathbf{x}) : \mathbb{R}^d \mapsto \mathbb{R}$ is modeled as a zero mean Gaussian process for simplicity; and the noise e is assumed to be Gaussian distributed with zero mean and variance σ_e^2 . Moreover, the noise terms at different data points are assumed to be mutually independent. The set of all unknown GP hyper-parameters is denoted by $\boldsymbol{\theta} \triangleq [\boldsymbol{\theta}_h^T, \sigma_e^2]^T$ and the dimension of $\boldsymbol{\theta}$ is assumed to be p .

Given a training data set $\mathcal{D} \triangleq \{\mathbf{X}, \mathbf{y}\}$, where $\mathbf{y} = [y_1, y_2, \dots, y_n]^T$ is the vector comprising the outputs and $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n]$ is the matrix comprising the input vectors, the aim is to compute the posterior distribution of $\mathbf{y}_* = [y_{*,1}, y_{*,2}, \dots, y_{*,n_*}]^T$ given the corresponding test inputs $\mathbf{X}_* = [\mathbf{x}_{*,1}, \mathbf{x}_{*,2}, \dots, \mathbf{x}_{*,n_*}]$. Here, we let $\mathcal{D}_* \triangleq \{\mathbf{X}_*, \mathbf{y}_*\}$ be the test data set. According to the definition of Gaussian processes given before, the joint prior distribution of the training output \mathbf{y} and test output \mathbf{y}_* can be written explicitly as:

$$\begin{bmatrix} \mathbf{y} \\ \mathbf{y}_* \end{bmatrix} \sim \mathcal{N} \left(\mathbf{0}, \begin{bmatrix} \mathbf{K}(\mathbf{X}, \mathbf{X}) + \sigma_e^2 \mathbf{I}_n, & \mathbf{K}(\mathbf{X}, \mathbf{X}_*) \\ \mathbf{K}(\mathbf{X}_*, \mathbf{X}), & \mathbf{K}(\mathbf{X}_*, \mathbf{X}_*) + \sigma_e^2 \mathbf{I}_{n_*} \end{bmatrix} \right),$$

where $\mathbf{K}(\mathbf{X}, \mathbf{X})$ is an $n \times n$ matrix of covariances among the training inputs; $\mathbf{K}(\mathbf{X}, \mathbf{X}_*)$ is an $n \times n_*$ matrix of covariances between the training inputs and test inputs; $\mathbf{K}(\mathbf{X}_*, \mathbf{X}_*)$ is an $n_* \times n_*$ matrix of covariances among the test inputs. Here, we let $\mathbf{K}(\mathbf{X}, \mathbf{X})$ be a short term of $\mathbf{K}(\mathbf{X}, \mathbf{X}; \boldsymbol{\theta}_h)$ when the kernel hyper-parameters have been trained and the associated optimization process is not the spotlight.

Applying the results of conditional Gaussian distribution, we can easily derive the posterior distribution as

$$p(\mathbf{y}_* | \mathcal{D}, \mathbf{X}_*; \boldsymbol{\theta}_h) \sim \mathcal{N}(\bar{\mathbf{m}}, \bar{\mathbf{V}}), \quad (3)$$

where the posterior mean and posterior variance are respectively,

$$\bar{\mathbf{m}} = \mathbf{K}(\mathbf{X}_*, \mathbf{X}) [\mathbf{K}(\mathbf{X}, \mathbf{X}) + \sigma_e^2 \mathbf{I}_n]^{-1} \mathbf{y}, \quad (4)$$

$$\begin{aligned} \bar{\mathbf{V}} &= \mathbf{K}(\mathbf{X}_*, \mathbf{X}_*) + \sigma_e^2 \mathbf{I}_{n_*} \\ &\quad - \mathbf{K}(\mathbf{X}_*, \mathbf{X}) [\mathbf{K}(\mathbf{X}, \mathbf{X}) + \sigma_e^2 \mathbf{I}_n]^{-1} \mathbf{K}(\mathbf{X}, \mathbf{X}_*). \end{aligned} \quad (5)$$

In general, temporal Gaussian processes take training input $\mathbf{x}_t = [x_{1,t}, x_{2,t}, \dots, x_{d,t}]^T$ with discrete time index $t = 1, 2, \dots, n$, where $x_{1,t}, x_{2,t}, \dots, x_{d,t}$ are specifically the d features observed at time t . In this paper, we focus on the one-dimensional (1-D) time series with $d = 1$ and $\mathbf{x}_t = x_t = t$.

B. Classic GP Hyper-parameter Optimization

Next, we introduce the classic ML based GP hyper-parameter estimation. Due to the Gaussian assumption on the noise, the log-likelihood function can be obtained in closed form. The GP hyper-parameters can be tuned equivalently by minimizing the negative log-likelihood function (ignoring the unrelated terms) as

$$\boldsymbol{\theta}_{ML} \triangleq \arg \min_{\boldsymbol{\theta}} l(\boldsymbol{\theta}) = \mathbf{y}^T \mathbf{C}^{-1}(\boldsymbol{\theta}) \mathbf{y} + \log \det(\mathbf{C}(\boldsymbol{\theta})), \quad (6)$$

where $\mathbf{C}(\boldsymbol{\theta}) \triangleq \mathbf{K}(\mathbf{X}, \mathbf{X}; \boldsymbol{\theta}_h) + \sigma_e^2 \mathbf{I}_n$. This optimization problem is mostly solved via gradient based methods, such as LFGS-Newton or conjugate gradient [2], which requires

the following partial derivatives for $i = 1, 2, \dots, p$ in closed form:

$$\frac{\partial l(\boldsymbol{\theta})}{\partial \theta_i} = \text{tr} \left(\mathbf{C}^{-1}(\boldsymbol{\theta}) \frac{\partial \mathbf{C}(\boldsymbol{\theta})}{\partial \theta_i} \right) - \mathbf{y}^T \mathbf{C}^{-1}(\boldsymbol{\theta}) \frac{\partial \mathbf{C}(\boldsymbol{\theta})}{\partial \theta_i} \mathbf{C}^{-1}(\boldsymbol{\theta}) \mathbf{y}.$$

C. Linear Multiple Kernel

Linear multiple kernel, as its name suggests, constitutes a linear combination of primitive kernels whose weights are to be optimized. In this paper, we solely focus on the scenario, in which the underlying kernel function $k(t, t')$ is **completely unknown** but is approximated as $k(t, t') \approx \sum_{i=1}^m \alpha_i k_i(t, t')$, where the basis sub-kernel functions $k_i(t, t')$, $i = 1, 2, \dots, m$ are known and the weights α_i , $i = 1, 2, \dots, m$ are the optimization variables, subject to $\alpha_i \geq 0$. Often, the number of the basis sub-kernels, m , is set large to allow for good approximation. For this scenario, no expert knowledge is required. The associated kernel hyper-parameters are $\boldsymbol{\theta}_h = \boldsymbol{\alpha} = [\alpha_1, \alpha_2, \dots, \alpha_m]^T$. We will introduce two ways of constructing a grid spectral mixture kernel in Section III with the aim to let the data decide on the most favorable stationary kernel function approximated by a linear multiple of basis kernels.

III. STATIONARY KERNEL DESIGN IN THE FREQUENCY DOMAIN

In subsection III-A, we first briefly review the spectral mixture (SM) kernel proposed originally in [18] for approximating any stationary kernel while stressing out the associated difficulties when optimizing with respect to the hyper-parameters. In subsection III-B, we introduce two ways of constructing grid spectral mixture (GSM) kernel for building 1-D temporal Gaussian process regression models. Lastly, we show how to combine Welch periodogram with L_1 norm regularization for advanced setup of the GSM kernel in subsection III-C.

A. SM Kernel [18]

The SM kernel undertakes approximation in the frequency domain using the fact that a stationary kernel function and its spectral density are Fourier duals due to the following corollary of Bochner's theorem given in [2].

Corollary 1. *For time series where the free variable is time, i.e., $\mathbf{x} = t$, $\tau = t - t'$, f being the normalized frequency (i.e., $f \in [0, 1/2)$) and in the case that the spectral density $S(f)$ exists, the stationary kernel function, $k(\tau)$, and its spectral density of the kernel function, $S(f)$, are Fourier duals of each other as shown below:*

$$k(\tau) = \int_{\mathbb{R}^1} S(f) \exp[j2\pi\tau f] df, \quad (7a)$$

$$S(f) = \int_{\mathbb{R}^1} k(\tau) \exp[-j2\pi\tau f] d\tau. \quad (7b)$$

The salient SM kernel is designed by approximating the spectral density, $S(f)$, of the underlying stationary kernel by a Gaussian mixture. Taking the inverse Fourier transform of $S(f)$, yields a stationary kernel in the time-domain as

$$k(t, t'; \boldsymbol{\theta}_h) = k(\tau) = \sum_{q=1}^Q \alpha_q \exp[-2\pi^2\tau^2\sigma_q^2] \cos(2\pi\tau\mu_q), \quad (8)$$

where $\theta_h \triangleq [\alpha_1, \dots, \alpha_Q, \mu_1, \dots, \mu_Q, \sigma_1^2, \dots, \sigma_Q^2]^T$ denotes the SM kernel hyper-parameters with Q being a fixed number of mixture components, and $\alpha_q, \mu_q, \sigma_q^2$ being the weight, mean and variance of the q -th mixture component, respectively. The SM kernel is able to approximate any stationary kernel arbitrarily well in L_1 norm according to the Wiener's theorem of approximation [23].

However, minimizing the negative log-likelihood with respect to θ in light of Eq.(6), it may easily get stuck at a bad local optimum, because the cost function is non-convex in terms of θ and may not have any favorable structure to facilitate the optimization process.

B. Proposed GSM Kernel

To address the potential numerical problems with the original SM kernel, we proposed a GSM kernel in [1] with the goal to modify the original SM kernel by fixing the μ and σ parameters to *a priori* selected values in a grid. To be precise, the spectral density is approximated by the GSM kernel as

$$S(f) = \sum_{i=1}^m \alpha_i s_i(f), \quad (9)$$

where each $s_i(f) = \mathcal{N}(f; \mu_i, \sigma_i^2) + \mathcal{N}(f; -\mu_i, \sigma_i^2)$ is evaluated at a fixed point in a grid (μ_i, σ_i^2) , sampled either uniformly or randomly from a two dimensional space confined in $[\mu_{low}, \mu_{high}]$ and $[\sigma_{low}^2, \sigma_{high}^2]$. The sampling strategies are shown in Fig. 1 for clarity. Taking the inverse Fourier transform of the above spectral density, $S(f)$, yields our first GSM kernel formulation in the time domain as

$$\begin{aligned} k(t, t'; \alpha) &= \sum_{i=1}^m \alpha_i k_i(t, t') \\ &= \sum_{i=1}^m \alpha_i \exp[-2\pi^2 \tau^2 \sigma_i^2] \cos(2\pi \tau \mu_i), \end{aligned} \quad (10)$$

where α is a vector of the kernel hyper-parameters, namely the unknown non-negative weights. Since the grids are generated in the 2-D (μ, σ) space, the resultant kernel is called 2-D GSM kernel.

The problem with the so designed 2-D GSM kernel lies in the large number of unknown parameters to be optimized. We note that the weights $\alpha_1, \alpha_2, \dots, \alpha_m$ are all non-negative numbers, but we will not constrain the sum $\sum_{i=1}^m \alpha_i$ to be equal to one for approximating a spectral density whose integral is not equal to one, as in [19]. Therefore, we slightly abuse the term ‘‘Gaussian mixture’’ mostly used for probability density approximation [24].

In the following, we derive a modified 1-D GSM kernel by fixing the variance parameters, σ_i , in Eq.(10) to a small fixed value, σ , so as to reduce the high model complexity. The resultant GSM kernel boils down to

$$k(t, t'; \alpha) = \sum_{i=1}^m \alpha_i \exp(-2\pi^2 \tau^2 \sigma^2) \cos(2\pi \tau \mu_i). \quad (11)$$

To differentiate with the 2-D GSM kernel, the kernel given in Eq.(11) is called 1-D GSM kernel, because the grids are generated in the 1-D μ -space, given a fixed σ . With this

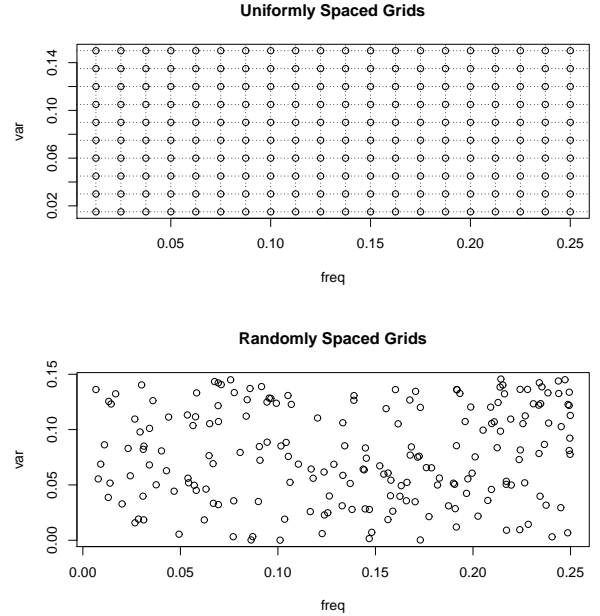


Fig. 1. Illustration of the two strategies for generating grids. In this specific example, μ_{low} is set to be 0, $\mu_{high} = 0.25$, $\sigma_{low} = 0$ and $\sigma_{high} = 0.15$.

1-D GSM kernel, the underlying spectral density, $S(f)$, is approximated by a linear weighted sum of Gaussian basis functions with varying shifts, μ_i , while fixed bandwidth, σ . In addition to the kernel design, we also demonstrate some useful properties of the proposed GSM kernels.

Theorem 1. *Some properties of the GSM kernel in Eq. (11) are given as follows:*

- 1) *It is a valid kernel.*
- 2) *It is smooth with derivatives of all orders.*
- 3) *Each one of the sub-kernel functions, $k_i(\tau)$, is square integrable for any $i = 1, 2, \dots, m$.*
- 4) *For big data set with size $n \gg \frac{4}{\pi\sigma}$, the sub-kernel matrix is sparse and close to a band matrix with equal lower and upper bandwidths (irrespective of μ_i), which enables more efficient utilization of computer memory, e.g., in MATLAB [25].*
- 5) *For a given data set with n samples, when the variance parameter, σ , is chosen sufficiently small, then for any frequency parameter $\mu_i \in [0, 1/2)$, the corresponding sub-kernel matrix has low rank, $\text{rank}(\mathbf{K}_i) \ll n$.*

Proof. Sketch of the proofs are summarized below:

- Proof of property (1) is given in the Appendix A.
- Verification of property (2) is straightforward.
- Reasoning of properties (3) and (4) is given in the supplement.
- Proof of property (5) is given in the Appendix B. □

It is easy to see that the above properties hold for any sub-kernel of the 2-D GSM kernel in Eq.(10) as well.

Remark 1. *Our way of constructing the 1-D GSM kernel is related, in some sense, to the non-parametric kernel density*

estimator using an optimal kernel width [26]. The difference, however, lies in the distribution of the “frequency variables”. For the 1-D GSM kernel, the frequency variables μ_i are selected either uniformly or randomly from the selected region; while in the non-parametric kernel density estimation, the “frequency variables” are essentially generated from the underlying density function to be reconstructed.

C. Advanced Setup of the 1-D GSM Kernel

In the previous subsection, we have seen a modified GSM kernel function as given in Eq. (11). In order to make it attractive from a practical point of view, as it will be confirmed in Section VI, one simply needs to 1) choose a moderate number of modes, m ; 2) set a small number σ common to all grids, and 3) sample μ_i , $i = 1, 2, \dots, m$, either uniformly or randomly from $[0, 1/2)$. Naturally, one may ask for more advanced setup of the GSM kernel with reduced model complexity, promising sampling areas, and better initial guess of the unknown weights.

For the purpose of obtaining an advanced setup, we could exploit the observations $y(t)$, $t = 1, 2, \dots, n$, which are assumed to comprise a noisy realization of the underlying stationary random process $f(t)$, so that to build an estimate of the true spectral density. A candidate is to use the Welch periodogram [27] as an estimator of the underlying spectral density, $S(f)$. To construct a Welch periodogram, we need to partition $y(t)$, $t = 1, 2, \dots, n$, into L overlapped segments, $y_l(t)$, $l = 1, 2, \dots, L$, each with only D data points. For each segment, a local periodogram is then computed as

$$I_{W,l}^D(f) = \frac{1}{DA} \left| \sum_{t=1}^D w(t) y_l(t) e^{-j2\pi ft} \right|^2, \quad (12)$$

where $w(t)$ is a deterministic window function, e.g., the Bartlett window, and $A = \frac{1}{D} \sum_{t=1}^D |w(t)|^2$ is a normalization factor. The final Welch periodogram, $\hat{S}_W(f)$, is given as the average of the L local periodograms. For a big data set with n , L , and D all being large, the Welch periodogram is an asymptotically consistent estimator of the underlying power spectral density. Hence, by inspecting the Welch periodogram, we may obtain good prior knowledge about the model complexity, m , as well as the salient areas for the sampling points in the grid.

Next, the previously obtained periodogram will be used to compute a potentially good initial guess of the weights, α . To this end, we solve

$$\min_{\alpha} \|s_W - \Psi\alpha\|_2^2 + \lambda \|\alpha\|_1, \quad (13)$$

where $s_W = [\hat{S}_W(\mu_1), \hat{S}_W(\mu_2), \dots, \hat{S}_W(\mu_m)]^T$ contains the periodogram values evaluated at the discrete frequencies, $\mu_1, \mu_2, \dots, \mu_m$; matrix Ψ is of size $m \times m$, whose i -th row is the transpose of $\mathbf{s}(\mu_i) = [s_1(\mu_i), s_2(\mu_i), \dots, s_m(\mu_i)]^T$ with each entry computed according to the definition of the Gaussian mixture component introduced in Eq. (9). A practical and efficient method for solving a large scale $L1$ -regularized least-squares problem in Eq.(13) can be found in [28].

Here, we must acknowledge that using empirical periodogram for additional information concerning the underlying

spectral density was already mentioned in [18]. However, in our current context, it is used as a potentially better initial guess of α , given the knowledge that our hyper-parameter estimate will be sparse.

IV. MEMORY EFFICIENT KERNEL MATRIX APPROXIMATIONS

When the proposed GSM kernel contains a large number of sub-kernels, i.e., m is large, and moreover the data size n is large, unaffordable memory is needed to store the m huge sub-kernel matrices during the hyper-parameter optimization process, as will be introduced in Section V. Often, a factor L_i , satisfying $\mathbf{K}_i = L_i L_i^T$, is stored instead of the sub-kernel matrix \mathbf{K}_i with much reduced memory, especially when \mathbf{K}_i has low rank. In this section, we discuss two kernel matrix approximations, namely the Nyström approximation in subsection IV-A and the random Fourier feature approximation in subsection IV-B, that can be adopted to provide good approximations of L_i with relatively low computational complexity and reduced memory.

A. Nyström Approximation [29]

First, we introduce the Nyström approximation [29] of the kernel matrix \mathbf{K}_i , $i = 1, 2, \dots, m$. In the sequel, we omit the subscript i for brevity because the same procedure can be applied to any \mathbf{K}_i . Detailed steps are as follows:

Step 1: Sample a subset of p ($\leq n$) training inputs to form $\tilde{\mathbf{X}}$ from the complete set of training inputs \mathbf{X} .

Step 2: Compute $\mathbf{K}^{(p)}$ with the sub-sampled training inputs $\tilde{\mathbf{X}}$. Herein, the superscript (p) indicates $\mathbf{K}^{(p)}$ is of size $p \times p$.

Step 3: Perform eigendecomposition of the smaller kernel matrix $\mathbf{K}^{(p)}$ as

$$\mathbf{K}^{(p)} = \sum_{l=1}^{\tilde{p}} \lambda_l^{(p)} \mathbf{u}_l^{(p)} \left(\mathbf{u}_l^{(p)} \right)^T, \quad (14)$$

where \tilde{p} denotes the effective number of eigenvalues that are distinctly larger than zero and obviously $\tilde{p} \leq p$. We further define $\Sigma^{(p)} \triangleq \text{diag}(\lambda_1^{(p)}, \lambda_2^{(p)}, \dots, \lambda_{\tilde{p}}^{(p)})$ and $\mathbf{U}^{(p)} \triangleq [\mathbf{u}_1^{(p)}, \mathbf{u}_2^{(p)}, \dots, \mathbf{u}_{\tilde{p}}^{(p)}]$ for later use.

Step 4: Apply Nyström approximation to the eigenvalues and eigenvectors obtained in the previous step as follows:

$$\tilde{\lambda}_l = \frac{n}{p} \lambda_l^{(p)}, \quad l = 1, 2, \dots, \tilde{p}, \quad (15)$$

$$\tilde{\mathbf{u}}_l = \sqrt{\frac{p}{n} \frac{1}{\lambda_l^{(p)}}} \mathbf{K}(\mathbf{X}, \tilde{\mathbf{X}}) \mathbf{u}_l^{(p)}, \quad l = 1, 2, \dots, \tilde{p}, \quad (16)$$

where $\tilde{\lambda}_l$ and $\tilde{\mathbf{u}}_l$ are respectively the approximated l -th eigenvalue and eigenvector of the original $n \times n$ kernel matrix \mathbf{K} , and $\mathbf{K}(\mathbf{X}, \tilde{\mathbf{X}})$ is an $n \times p$ matrix of correlations between the training inputs \mathbf{X} and sub-sampled training inputs $\tilde{\mathbf{X}}$.

Step 5: Finally, we obtain a low-rank (of rank \tilde{p}) approximation of the original kernel matrix \mathbf{K} as follows:

$$\mathbf{K} \approx \tilde{\mathbf{K}} = \sum_{l=1}^{\tilde{p}} \tilde{\lambda}_l \tilde{\mathbf{u}}_l \tilde{\mathbf{u}}_l^T = \tilde{\mathbf{U}} \tilde{\Sigma} \tilde{\mathbf{U}}^T, \quad (17)$$

where $\tilde{\mathbf{U}} \triangleq [\tilde{\mathbf{u}}_1, \tilde{\mathbf{u}}_2, \dots, \tilde{\mathbf{u}}_{\tilde{p}}]$ is the matrix of the \tilde{p} eigenvectors and $\tilde{\mathbf{\Sigma}} \triangleq \text{diag}(\tilde{\lambda}_1, \tilde{\lambda}_2, \dots, \tilde{\lambda}_{\tilde{p}})$ is a diagonal matrix of the \tilde{p} eigenvalues. Lastly, we approximate the factor \mathbf{L} by $\tilde{\mathbf{L}} \triangleq \tilde{\mathbf{U}}\tilde{\mathbf{\Sigma}}^{1/2}$, which is of smaller size $n \times \tilde{p}$.

It is easy to verify that the memory usage for storing $\tilde{\mathbf{L}}$ is reduced to $\tilde{p}/n \times 100\%$ of the original usage for storing \mathbf{L} . Moreover, the computational complexity for performing eigendecomposition is also reduced from $\mathcal{O}(n^3)$ to $\mathcal{O}(\tilde{p}^3)$.

B. Random Fourier Feature Approximation [30]

Next, we introduce the random Fourier feature approximation. When the spectral density function $S(f)$ is an even function of f , we can easily derive the corresponding stationary kernel function as

$$k(t, t') = \mathbb{E}_{S(f)} [\cos(2\pi ft - 2\pi ft')]. \quad (18)$$

By replacing the above integration with Monte-Carlo integration, $k(t, t')$ is approximated as follows:

$$k(t, t') \approx \frac{1}{R} \sum_{r=1}^R \cos(2\pi f_r t - 2\pi f_r t'), \quad (19)$$

where $f_r, r = 1, 2, \dots, R$ are sampled from the spectral density function $S(f)$. Let $\omega_r \triangleq 2\pi f_r, r = 1, 2, \dots, R$ and define

$$\phi_\omega(t) \triangleq \frac{1}{\sqrt{R}} [\cos(\omega_1 t), \sin(\omega_1 t), \dots, \cos(\omega_R t), \sin(\omega_R t)]^T \quad (20)$$

we have $k(t, t') \approx \phi_\omega^T(t) \phi_\omega(t')$.

The GSM kernel proposed in the above subsection is of form $k(t, t'; \boldsymbol{\alpha}) = \sum_{i=1}^m \alpha_i k_i(t, t')$ and it can be approximated as

$$k(t, t'; \boldsymbol{\alpha}) \approx \sum_{i=1}^m \alpha_i \phi_{i,\omega}^T(t) \phi_{i,\omega}(t'), \quad (21)$$

by applying random Fourier representation to each sub-kernel function, i.e.,

$$k_i(t, t') = \exp[-2\pi^2 \tau^2 \sigma_i^2] \cos(2\pi \tau \mu_i) \approx \phi_{i,\omega}^T(t) \phi_{i,\omega}(t'), \quad (22)$$

where the random Fourier features for the i -th sub-kernel $k_i(t, t')$ are randomly sampled from $s_i(f)$ (also a valid distribution function).

With the aid of the random Fourier feature representation of the sub-kernel, the overall GSM kernel matrix can be represented as $\mathbf{K} = \sum_{i=1}^m \alpha_i \mathbf{K}_i \approx \tilde{\mathbf{K}} = \sum_{i=1}^m \alpha_i \tilde{\mathbf{L}}_i \tilde{\mathbf{L}}_i^T$, where $\tilde{\mathbf{L}}_i = [\phi_{i,\omega}(t_1), \phi_{i,\omega}(t_2), \dots, \phi_{i,\omega}(t_n)]^T$ is an $n \times 2R$ matrix and $\mathbf{L}_i \approx \tilde{\mathbf{L}}_i$. The memory usage for storing $\tilde{\mathbf{L}}$ is reduced to $2R/n \times 100\%$ of the original usage for storing \mathbf{L} . The computational complexity is mainly due to a batch of samplings from a Gaussian distribution, which remains low. Random Fourier feature is widely used for resource limited kernel approximations, e.g., for fast online learning [31] and others [6]. Alternative to the random Fourier features, one may also use the Fastfood features that can be computed more efficiently [32].

C. RAE versus Storage

To compare the above two methods in terms of approximation accuracy and storage, we adopt the widely used metric relative approximation error (RAE) as the performance metric, which is given by $\|\mathbf{K} - \tilde{\mathbf{K}}\|_F / \|\mathbf{K}\|_F$, where \mathbf{K} is the exact kernel matrix and $\tilde{\mathbf{K}} = \tilde{\mathbf{L}}\tilde{\mathbf{L}}^T$ is its approximation. Due to space limitation, the results are shown in the supplement.

The general conclusions are as follows:

- For small data set, Nyström approximation may require less memory than the random Fourier feature approximation in order to achieve a similar small value of RAE, say less than 1%.
- For medium or large data set, random Fourier feature approximation may require less memory than the Nyström approximation in order to achieve a similar small value of RAE. This is because the number of random features needed for constructing a good approximation can be kept to several hundreds, not sensitive to sample size of the selected data set; while the number of data points needed by the Nyström approximation increases with the sample size, in general. When the kernel matrices have low rank, less samples is needed by the Nyström approximation.

V. ML BASED GP HYPER-PARAMETER OPTIMIZATION

For linear multiple kernel, including the proposed GSM kernel, the associated maximum-likelihood based GP hyper-parameter optimization problem can be cast as

$$\boldsymbol{\theta}_{ML} = \arg \min_{\boldsymbol{\alpha}, \sigma_e^2} \{ \mathbf{y}^T \mathbf{C}^{-1}(\boldsymbol{\alpha}, \sigma_e^2) \mathbf{y} + \log \det(\mathbf{C}(\boldsymbol{\alpha}, \sigma_e^2)) \}, \quad (23)$$

subject to $\boldsymbol{\alpha} \geq \mathbf{0}$ and $\sigma_e^2 \geq 0$. Here, $\boldsymbol{\theta} = [\boldsymbol{\alpha}^T, \sigma_e^2]^T$ and $\mathbf{C}(\boldsymbol{\alpha}, \sigma_e^2) \triangleq \sum_{i=1}^m \alpha_i \mathbf{K}_i + \sigma_e^2 \mathbf{I}_n$, where \mathbf{K}_i is the i -th sub-kernel matrix evaluated at the grid point (μ_i, σ_i) for the 2-D GSM kernel or (μ_i, σ) for the 1-D GSM kernel. The cost function in Eq.(23) is a difference of two convex functions with respect to $\boldsymbol{\alpha}$ and σ_e^2 , therefore the optimization problem belongs to the well known difference-of-convex program (DCP) [33]–[35]. Here, we want to stress, once more, that the primary idea behind the newly proposed GSM kernel is to maintain good approximation capability with a structure that leads to a well known optimization problem with respect to the GP hyper-parameters. However, ML method for the previously suggested SM kernel leads to a general non-convex hyper-parameter optimization task. The additional structure in Eq.(23) can facilitate the optimization task in terms of convergence speed and avoidance of bad local minimum, as will be seen in our experiments.

In the following, we derive two numerical methods for optimizing the GP hyper-parameters. In subsection V-A, we derive a sequential majorization-minimization (MM) method. In subsection V-B, we derive a nonlinearly constrained alternating direction method of multipliers (ADMM) [36]. No matter which method is adopted, the solution can be proven to be sparse according to the theorem provided along with some other properties in subsection V-C.

Algorithm 1: Sequential MM Method

Input: \mathbf{y} and $\mathbf{L}_i, i = 1, 2, \dots, m$
Output: $\boldsymbol{\theta}_{ML}$

- 1 Initialization: $k = 0, \boldsymbol{\theta}^0$
 - 2 **while** the convergence condition is not satisfied **do**
 - 3 Compute $h(\boldsymbol{\theta}^k)$ and the gradient $\nabla_{\boldsymbol{\theta}} h(\boldsymbol{\theta}^k)$.
 - 4 Solve Eq.(25) for $\boldsymbol{\theta}^{k+1}$.
 - 5 Set $k = k + 1$.
 - 6 **end**
 - 7 $\boldsymbol{\theta}_{ML} = \boldsymbol{\theta}^k$
-

A. Sequential MM Method [34]

The main idea of the MM method is to solve $\min_{\boldsymbol{\theta} \in \Theta} l(\boldsymbol{\theta})$ with $\Theta \subseteq \mathbb{R}^{m+1}$ through an iterative scheme, where at each iteration a so-called majorization function $\bar{l}(\boldsymbol{\theta}, \boldsymbol{\theta}^k)$ of $l(\boldsymbol{\theta})$ at $\boldsymbol{\theta}^k \in \Theta$ is minimized, i.e.,

$$\boldsymbol{\theta}^{k+1} = \arg \min_{\boldsymbol{\theta} \in \Theta} \bar{l}(\boldsymbol{\theta}, \boldsymbol{\theta}^k), \quad (24)$$

where $\bar{l}: \Theta \times \Theta \rightarrow \mathbb{R}$ satisfies $\bar{l}(\mathbf{x}, \mathbf{x}) = l(\mathbf{x})$ for $\mathbf{x} \in \Theta$ and $l(\mathbf{x}) \leq \bar{l}(\mathbf{x}, \mathbf{z})$ for $\mathbf{x}, \mathbf{z} \in \Theta$. For this particular DCP problem in Eq.(23), $l(\boldsymbol{\theta}) = g(\boldsymbol{\theta}) - h(\boldsymbol{\theta})$, where $g(\boldsymbol{\theta}) = \mathbf{y}^T \mathbf{C}(\boldsymbol{\theta})^{-1} \mathbf{y}$, $h(\boldsymbol{\theta}) = -\log \det \mathbf{C}(\boldsymbol{\theta})$ and $\mathbf{C}(\boldsymbol{\theta}) = \sum_{i=1}^m \alpha_i \mathbf{L}_i \mathbf{L}_i^T + \sigma_e^2 \mathbf{I}_n$; $g(\boldsymbol{\theta}), h(\boldsymbol{\theta}): \Theta \rightarrow \mathbb{R}$ are convex and differentiable functions with Θ being a convex set in \mathbb{R}^{m+1} . Here, we use the so-called linear majorization, i.e., we make the convex function $h(\boldsymbol{\theta})$ affine by performing the first-order Taylor expansion and obtain $\bar{l}(\boldsymbol{\theta}, \boldsymbol{\theta}^k) = g(\boldsymbol{\theta}) - h(\boldsymbol{\theta}^k) - \nabla_{\boldsymbol{\theta}}^T h(\boldsymbol{\theta}^k)(\boldsymbol{\theta} - \boldsymbol{\theta}^k)$. Hence, at each iteration, minimizing the cost function in Eq.(24) becomes a convex optimization problem. The MM method is guaranteed to converge to a stationary point when some regularization conditions are satisfied [34]. But it is noticed that solving this problem directly with CVX, a package for specifying and solving convex programs [37], is very computationally demanding.

Since $g(\boldsymbol{\theta})$ is a matrix fractional function, each iteration in the MM method actually solves a convex matrix fractional minimization problem. Due to the fact that $\sum_{i=1}^m \alpha_i \mathbf{L}_i \mathbf{L}_i^T + \sigma_e^2 \mathbf{I}_n$ is a sum of positive semi-definite terms, the SDP problem can be further cast as a conic quadratic optimization problem with $m + 1$ rotated quadratic cone constraints, i.e.,

$$\begin{aligned} & \min_{\mathbf{z}, \boldsymbol{\theta}, \mathbf{v}, \mathbf{w}} \quad 2(\mathbf{1}^T \mathbf{z}) - \nabla_{\boldsymbol{\theta}}^T h(\boldsymbol{\theta}^k) \boldsymbol{\theta} \\ \text{s.t.} \quad & \|\mathbf{w}_i\|_2^2 \leq 2\theta_i z_i, \quad i = 1, 2, \dots, m \\ & \|\mathbf{v}\|_2^2 \leq 2z' \\ & \mathbf{y} = \sum_{i=1}^m \mathbf{L}_i \mathbf{w}_i + \sigma_e \mathbf{v}, \quad \boldsymbol{\theta} \geq \mathbf{0}, \mathbf{z} \geq \mathbf{0} \end{aligned} \quad (25)$$

where $\boldsymbol{\theta} \in \mathbb{R}^{m+1}$, $\mathbf{z} = [z_1, z_2, \dots, z_m, z']^T \in \mathbb{R}^{m+1}$, $\mathbf{v} \in \mathbb{R}^n$, and $\mathbf{w}_i \in \mathbb{R}^{n_i}$ for $i = 1, 2, \dots, m$. The conic quadratic optimization problem here is equivalent to a second-order cone program that can be solved efficiently using the commercial solver MOSEK [34].

Remark 2. The computational complexity for solving one iteration of the above second-order cone program scales as

$\mathcal{O}(n^2 \cdot \max(n, \sum_{i=1}^m n_i))$, where n_i stands for the rank of \mathbf{K}_i . The worst case complexity is $\mathcal{O}(mn^3)$ if all GSM sub-kernel matrices have full rank. Fortunately, as it was reported in [34] the MM method requires only a few iterations to achieve a good local optimum in practice.

Remark 3. The above MM method matches perfectly with the proposed GSM kernel. This is due to the fourth property of the GSM kernel given in Theorem 1, i.e., for a given number of data samples, n , and a sufficiently small σ , the rank of \mathbf{L}_i satisfies $n_i \ll n$, making $\sum_{i=1}^m n_i$ relatively small. Moreover, the two matrix approximation approaches are also helpful for reducing the complexity. For instance, using the random Fourier feature approximation can reduce the computational complexity to $\mathcal{O}(n^2 \cdot \max(n, 2mR))$, where R is the number of random features, specified in Section IV.

B. Nonlinearly Constrained ADMM

In this subsection, we will propose a nonlinearly constrained ADMM for solving the optimal GP hyper-parameters, $\boldsymbol{\theta}$, from the maximum-likelihood estimation problem in Eq. (23). This new method has good potential to find a better local minimum with smaller negative likelihood value, $l(\boldsymbol{\theta})$, and the prediction MSE as compared to the sequential MM method and the classic gradient descent method. However, this method constrains itself to time series with short data records because its sub-problems involve matrix inversion and matrix multiplications, which scale as $\mathcal{O}(n^3)$ in general.

The idea is as follows. We reformulate the original problem by introducing an $n \times n$ matrix \mathbf{S} and solve instead

$$\arg \min_{\mathbf{S}, \boldsymbol{\alpha}} \mathbf{y}^T \mathbf{S} \mathbf{y} - \log \det(\mathbf{S}), \quad (26)$$

subject to $\mathbf{S} (\sum_{i=1}^m \alpha_i \mathbf{K}_i + \sigma_e^2 \mathbf{I}_n) = \mathbf{I}_n$ and $\boldsymbol{\alpha} \geq \mathbf{0}$. Although $\sigma_e^2 \geq 0$ can be estimated jointly, we simply assume it is known *a priori* and focus on the kernel hyper-parameters, $\boldsymbol{\alpha}$. This is for ease of notation and narration in the sequel.

The augmented Lagrangian function is then formulated as:

$$\begin{aligned} L_{\rho}(\mathbf{S}, \boldsymbol{\alpha}, \boldsymbol{\Lambda}) &= \mathbf{y}^T \mathbf{S} \mathbf{y} - \log \det(\mathbf{S}) \\ &+ \left\langle \boldsymbol{\Lambda}, \mathbf{S} \left(\sum_i \alpha_i \mathbf{K}_i + \sigma_e^2 \mathbf{I}_n \right) - \mathbf{I}_n \right\rangle \\ &+ \frac{\rho}{2} \left\| \mathbf{S} \left(\sum_i \alpha_i \mathbf{K}_i + \sigma_e^2 \mathbf{I}_n \right) - \mathbf{I}_n \right\|_F^2, \end{aligned} \quad (27)$$

where the regularization parameter $\rho > 0$ is fixed *a priori*. The ADMM applied to Eq.(27) iteratively decomposes into solving the following sub-problems:

$$\mathbf{S}^{k+1} = \arg \min_{\mathbf{S}} L_{\rho}(\mathbf{S}, \boldsymbol{\alpha}^k, \boldsymbol{\Lambda}^k) \quad (28)$$

$$\alpha_i^{k+1} = \arg \min_{\alpha_i} L_{\rho}(\mathbf{S}^{k+1}, \{\alpha_i, \boldsymbol{\alpha}_{-i}^{k,k+1}\}, \boldsymbol{\Lambda}^k), i = 1, \dots, m \quad (29)$$

$$\boldsymbol{\Lambda}^{k+1} = \boldsymbol{\Lambda}^k + \rho' \left[\mathbf{S}^{k+1} \left(\sum_i \alpha_i^{k+1} \mathbf{K}_i + \sigma_e^2 \mathbf{I}_n \right) - \mathbf{I}_n \right], \quad (30)$$

where $\boldsymbol{\alpha}_{-i}^{k,k+1} \triangleq [\alpha_1^{k+1}, \alpha_2^{k+1}, \dots, \alpha_{i-1}^{k+1}, \alpha_{i+1}^k, \dots, \alpha_m^k]^T$ in Eq.(29).

Remark 4. It is not difficult to verify that the subproblems of the proposed ADMM in Eq. (28) and Eq. (29) are both convex in terms of the corresponding optimization variables.

Remark 5. Different from the conventional ADMM, in the shown nonlinearly constrained ADMM, ρ' used for the dual variable update in Eq.(30) is chosen to be smaller than ρ used for the primal update in Eq.(27). This novel configuration was first applied in the flexible proximal ADMM for consensus problems in [38], where the authors set $\rho = \rho' + L$ with L being the Lipschitz constant of the the gradient of the objective function and harvested improved convergence performance.

In order to avoid the high computational cost for solving \mathbf{S}^{k+1} precisely from Eq.(28), which involves solving a quadratic matrix equation, we resort to the steepest descent method, which is computationally cheaper. We numerically update

$$\mathbf{S}^{k+1,\eta+1} = \mathbf{S}^{k+1,\eta} + \mu^\eta d^\eta, \quad \eta = 0, 1, \dots, It_S - 1, \quad (31)$$

where $d^\eta = -\frac{\nabla_{\mathbf{S}} L_\rho}{\|\nabla_{\mathbf{S}} L_\rho\|_F}$, $\mathbf{S}^{k+1,0} := \mathbf{S}^k$, and It_S is a fixed number of inner iterations. The gradient of $L_\rho(\mathbf{S}, \boldsymbol{\alpha}^k, \boldsymbol{\Lambda}^k)$ with respect to \mathbf{S} , short as $\nabla_{\mathbf{S}} L_\rho$, is equal to

$$\begin{aligned} \nabla_{\mathbf{S}} L_\rho(\mathbf{S}, \boldsymbol{\alpha}^k, \boldsymbol{\Lambda}^k) &= 2\mathbf{y}\mathbf{y}^T - (\mathbf{y}\mathbf{y}^T) \circ \mathbf{I}_n - 2\mathbf{S}^{-1} + \mathbf{S}^{-1} \circ \mathbf{I}_n \\ &\quad + \boldsymbol{\Lambda}^k \mathbf{C}^k + (\boldsymbol{\Lambda}^k \mathbf{C}^k)^T - (\boldsymbol{\Lambda}^k \mathbf{C}^k) \circ \mathbf{I}_n \\ &\quad + \rho(\mathbf{S} \mathbf{C}^k \mathbf{C}^k + \mathbf{C}^k \mathbf{C}^k \mathbf{S} - (\mathbf{S} \mathbf{C}^k \mathbf{C}^k) \circ \mathbf{I}_n) \\ &\quad - \rho(2\mathbf{C}^k - \mathbf{C}^k \circ \mathbf{I}_n), \end{aligned} \quad (32)$$

where both $\mathbf{C}^k = \sum_i \alpha_i^k \mathbf{K}_i + \sigma_e^2 \mathbf{I}_n$ and \mathbf{S} are symmetric. The above gradient involves a matrix inverse of current \mathbf{S}^k which is computationally demanding. For speed up, we replace $(\mathbf{S}^k)^{-1}$ with \mathbf{C}^k as approximation in Eq.(32) whenever possible. To be precise, at each iteration, we stick to the approximated gradient if $\|\mathbf{S}^k \mathbf{C}^k - \mathbf{I}\|_F \leq \delta$, where δ is a manually selected threshold to trade-off approximation error and computational time; Otherwise, the original gradient in Eq. (32) will be used. The stepsize μ^η is selected according to Armijo rule [39] at each iteration. More details about the stepsize selection can be found in the supplement.

For solving α_i from Eq.(29), we take the derivative of $L_\rho(\mathbf{S}^{k+1}, \{\alpha_i, \boldsymbol{\alpha}_{-i}^{k,k+1}\}, \boldsymbol{\Lambda}^k)$ with respect to α_i , $\forall i = 1, 2, \dots, m$ and set it equal to zero, yielding

$$\begin{aligned} &\langle \boldsymbol{\Lambda}^k, \mathbf{S}^{k+1} \mathbf{K}_i \rangle + \rho [\alpha_i \cdot \text{tr}(\mathbf{K}_i^T \mathbf{S}^{k+1,T} \mathbf{S}^{k+1} \mathbf{K}_i)] \\ &+ \rho \cdot \text{tr} \left[\left(\tilde{\mathbf{K}}_{-i} \mathbf{K}_i + \sigma_e^2 \mathbf{K}_i \right) \mathbf{S}^{k+1,T} \mathbf{S}^{k+1} - \mathbf{S}^{k+1} \mathbf{K}_i \right] = 0. \end{aligned} \quad (33)$$

where $\tilde{\mathbf{K}}_{-i} = \sum_{j=1}^{i-1} \alpha_j^{k+1} \mathbf{K}_j + \sum_{j=i+1}^m \alpha_j^k \mathbf{K}_j$. Following the steps sketched in Appendix C, α_i^{k+1} can be re-expressed as

$$\alpha_i^{k+1} = \left[\alpha_i^k + \frac{\text{tr} \left[\mathbf{K}_i \mathbf{S}^{k+1} \left(\mathbf{I}_n - \mathbf{S}^{k+1} \tilde{\mathbf{C}}_i^{k+1} - \frac{1}{\rho} \boldsymbol{\Lambda}^k \right) \right]}{\text{tr}(\mathbf{K}_i \mathbf{S}^{k+1} \mathbf{S}^{k+1} \mathbf{K}_i)} \right]_+ \quad (34)$$

where

$$\tilde{\mathbf{C}}_i^{k+1} = \sum_{j=1}^{i-1} \alpha_j^{k+1} \mathbf{K}_j + \sum_{j=i}^m \alpha_j^k \mathbf{K}_j + \sigma_e^2 \mathbf{I}_n. \quad (35)$$

Algorithm 2: Proposed Nonlinearly Constrained ADMM

Input: \mathbf{y} , σ_e^2 , and $\mathbf{K}_i, i = 1, 2, \dots, m$

Output: $\boldsymbol{\alpha}_{ML}$

```

1 Initialization:  $k = 0, \boldsymbol{\alpha}^0, \boldsymbol{\Lambda}^0, \rho, \rho', \epsilon_{ADMM}, \epsilon_S, It_S$ .
2 Set  $\mathbf{C}^0 = \sum_i \alpha_i^0 \mathbf{K}_i + \sigma_e^2 \mathbf{I}_n, \mathbf{S}^0 = [\mathbf{C}^0]^{-1}$ 
3 for (outer iterations)  $k = 0, 1, \dots$  do
4    $\eta = 0, \mathbf{S}^{k+1,\eta=0} = \mathbf{S}^k$ 
5   for (inner iterations)  $\eta = 0, 1, \dots, It_S - 1$  do
6     1. Compute  $d^\eta = -\frac{\nabla_{\mathbf{S}} L_\rho}{\|\nabla_{\mathbf{S}} L_\rho\|_F}$  analytically
       according to Eq.(32) or its approximation
       obtained using  $\mathbf{C}^k$  to replace the inverse of  $\mathbf{S}^k$ .
7     2. Adopt Armijo rule to select the step size  $\mu^\eta$ 
       and perform:
8      $\mathbf{S}^{k+1,\eta+1} = \mathbf{S}^{k+1,\eta} + \mu^\eta d^\eta$ 
9     if  $\|\mathbf{S}^{k+1,\eta+1} - \mathbf{S}^{k+1,\eta}\|_F \leq \epsilon_S$  then
10       $\eta = \eta + 1$ 
11      break
12     end
13      $\eta = \eta + 1$ 
14   end
15   Update  $\mathbf{S}^{k+1} = \mathbf{S}^{k+1,\eta}$ 
16   for  $i = 1$  to  $m$  do
17     Compute  $\alpha_i^{k+1}$  analytically according to Eq. (34).
18     Compute  $\tilde{\mathbf{C}}_i^{k+1}$  analytically according to Eq. (35).
19   end
20   if  $\|\boldsymbol{\alpha}^{k+1} - \boldsymbol{\alpha}^k\| \leq \epsilon_{ADMM}$  then
21      $\boldsymbol{\alpha}_{ML} = \boldsymbol{\alpha}^{k+1}$ 
22     return
23   end
24   Update  $\mathbf{C}^{k+1} = \tilde{\mathbf{C}}_m^{k+1}$ 
25   Update  $\boldsymbol{\Lambda}^{k+1}$  analytically according to Eq. (30).
26   Set  $k = k + 1$ .
27 end
28  $\boldsymbol{\alpha}_{ML} = \boldsymbol{\alpha}^k$ 

```

It is noted that $\mathbf{K}_i^T = \mathbf{K}_i$ and $(\mathbf{S}^{k+1})^T = \mathbf{S}^{k+1}$ due to the symmetric property of a kernel matrix. For clarity, we provide detailed steps for implementing the proposed nonlinearly constrained ADMM in Algorithm 2.

Remark 6. When taking the initial guess $\boldsymbol{\Lambda}^0$ close to the optimal Lagrange multiplier $\boldsymbol{\Lambda}^*$ and taking ρ sufficiently large, solving the unconstrained minimization problem $L_\rho(\mathbf{S}, \boldsymbol{\alpha}, \boldsymbol{\Lambda})$ can yield points close to the local minimum \mathbf{S}^* and $\boldsymbol{\alpha}^*$ that satisfy the sufficient optimality conditions. Details can be found in sections 4.2 and 5.2 of [39].

C. Properties of the Optimized Hyper-Parameters

This subsection aims to give some additional properties of the optimized GP hyper-parameters from both the statistical signal processing and optimization perspectives.

Theorem 2. Global minimum $(\boldsymbol{\alpha}^*, \sigma_{e,*}^2)$ exists that leads Eq.(23) to minus infinity, when the output satisfies $\mathbf{y} = \mathbf{V}_i \mathbf{z}$ for $\mathbf{z} \in \mathbb{R}^p$, $\|\mathbf{z}\|_2^2 < \infty$, and $p < n$, where the $n \times p$ matrix $\mathbf{V}_i \triangleq \mathbf{U}_i \boldsymbol{\Sigma}_i^{1/2}$ with $\boldsymbol{\Sigma}_i^{1/2}$ being the diagonal matrix of

square-root of the p non-zero eigenvalues and \mathbf{U}_i of size $n \times p$ containing the corresponding eigenvectors of a rank-deficient sub-kernel matrix \mathbf{K}_i .

Proof. The proof can be found in [1] or in the supplement. \square

Theorem 3. *Every local minimum of Eq.(23) is achieved at a sparse solution, regardless of whether noise is present or not.*

Proof. See [40, Theorem 2]. \square

Remark 7. *The sparseness of the ML solution according to the above theorem is celebrating for the proposed 1-D GSM kernel. The reasons are twofold. First, it means that only the frequencies thought to be important for modeling by the data will be pinpointed, endowing good interpretation of the kernel. Second, by avoiding to use all the grids (or model freedom) to fit the data, over-parameterization problem as indicated by Proposition 1 can be effectively alleviated.*

VI. EXPERIMENTAL RESULTS

In this section, we aim to investigate the prediction performance of the proposed GSM kernel based GP and compare it with the SM kernel based GP proposed by Wilson *et.al.* in [18] and the sparse spectrum GP proposed by Lázaro-Gredilla *et.al.* in [16] from various aspects. We picked up in total 8 classic time series data sets for test. Descriptions of the data are shown in Table I. The training data, \mathcal{D} , is used for optimizing the GP hyper-parameters; while the test data, \mathcal{D}_* , is used for evaluating the prediction MSE.

TABLE I
DETAILS OF THE SELECTED DATA SETS.

| Name | Description | Training \mathcal{D} | Test \mathcal{D}_* |
|--------------|--|------------------------|----------------------|
| ECG | Electrocardiography of an ordinary person measured over a period of time | 680 | 20 |
| CO2 | CO2 concentration made between 1958 and the end of 2003 | 481 | 20 |
| Electricity | Monthly average residential electricity usage in Iowa City 1971-1979 | 86 | 20 |
| Employment | Wisconsin employment time series, trade, Jan. 1961 Oct. 1975 | 158 | 20 |
| Hotel | Monthly hotel occupied room average 1963-1976 | 148 | 20 |
| Passenger | Passenger miles (Mil) flown domestic U.K., Jul. 1962-May 1972 | 98 | 20 |
| Clay | Monthly production of clay bricks: million units. Jan 1956 Aug 1995 | 450 | 20 |
| Unemployment | Monthly U.S. female (16-19 years) unemployment figures (thousands) 1948-1981 | 380 | 20 |

A. Algorithmic Setup

For the proposed GSM kernel based GP, short for GSMGP, we provide its setup in each individual subsection. Source code and all test data are available online.¹

¹https://github.com/Paalis/MATLAB_GSM

The SM kernel based GP, short for SMGP, proposed by Wilson *et.al.*:

- We use the source code provided on the author’s web page and follow the default setup suggested therein.²
- We follow the initialization strategy given on the author’s web page as well. Random restart is, however, not used.
- The number of Gaussian mixture components Q is chosen to be 10 or 500 for the SM kernel.
- The SMGP model hyper-parameters are determined by a gradient-descent type method.

The Sparse spectrum (SS) GP, short for SSGP, proposed by Lázaro-Gredilla *et.al.*:

- We use the source code provided on the author’s web page and follow the default setup³.
- We follow the strategy given by the authors to initialize the hyper-parameters. The number of basis is set to $m = 500$ in the simulations.
- The SSGP model hyper-parameters are determined by a conjugate-gradient method.

It is noteworthy that the independent noise variance parameter σ_e^2 is estimated using the cross-validation filter type method [41] and it is kept common to all above GP models for fair comparisons. In the following experiments, we solely compare the performance of the GSM kernel, SM kernel, and SS kernel. In [16], [18], [19], extensive experiments with both synthesized and real data have confirmed the effectiveness of the SM kernel and SS kernel as compared to the elementary kernels such as the SE kernel and Matern kernel.

B. Performance of the 2-D GSM kernel with MM Method

This subsection is a wrap-up of the results obtained in [1] for the 2-D GSM kernel using a big number of grids generated from 2-D space. Therein, the test involved 30 independent Monte-Carlo (MC) runs, and in each MC run, a new set of 20,000 grid points were randomly generated in the 2-D $(\mu - \sigma)$ space confined by $\mu_{low} = 0$, $\mu_{high} = 0.5$, $\sigma_{low}^2 = 0$ and $\sigma_{high}^2 = 0.15$. We initialize the weights of the sub-kernels, α , to a vector of zeros for the 2-D GSM kernel. The prediction MSE is evaluated for all selected GP models. Besides, we count the number of MC runs (out of 30 in total), in which one method stucked at a bad/meaningless local minimum (i.e., does not provide a meaningful prediction) and calculate the ratio, referred to as program fail rate (PFR) in this paper. Note that, the meaningless results were excluded when we compute the MSE.

From the results shown in Table II, we can conclude that the proposed 2-D GSM kernel based GP regression has gained well improved prediction MSE and stability as compared to its competitors. We did not show the PFR of the SSGP because it can always get the trend/envelop of the data but fail to fit small-scaled, fine structures. Whereas, the SMGP using $Q = 10$ Gaussian modes can better fit the data with a good starting point but it may even fail to capture the trend of the data with a bad starting point. The performance of the

²<https://people.orie.cornell.edu/andrew/code/>

³<http://www.tsc.uc3m.es/~miguel/downloads.php>

TABLE II

PERFORMANCE COMPARISON BETWEEN THE PROPOSED GSMGP (WITH 2-D GRIDS) AND ITS COMPETITORS, SSGP AND SMGP, IN TERMS OF THE MSE AND THE PFR.

| Name | SSGP MSE | SMGP MSE | SMGP PFR | GSMGP MSE | GSMGP PFR |
|-------------|----------|----------|----------|-----------|-----------|
| ECG | 1.6E-01 | 2.1E+00 | 0.63 | NA | NA |
| CO2 | 2.0E+02 | 7.4E+04 | 0.83 | NA | NA |
| Electricity | 8.2E+03 | 1.8E+04 | 0.47 | 6.8E+03 | 0.2 |
| Employment | 7.7E+01 | 2.3E+04 | 0.27 | 3.9E+01 | 0.07 |
| Hotel | 1.9E+04 | 2.6E+05 | 0.33 | 2.4E+03 | 0 |
| Passenger | 6.9E+02 | 3.5E+03 | 0.37 | 1.7E+02 | 0 |
| Clay | 5.3E+02 | 4.8E+03 | 0.93 | NA | NA |
| Unemploy | 2.1E+04 | 1.2E+05 | 0.9 | NA | NA |

proposed GP model becomes better and more stable, when the number of the grids grows beyond around 10,000. In Table II, the results of the GSMGP on the *CO2*, *clay*, and *unemployment* data sets are not available since the large size of the unknown weights, $\dim(\alpha)$, and long data record jointly make the program beyond the processing capability of our computer.⁴ Apart from the improved performance, the average number of non-zero α values generated by the ML method is equal to 26, 19, 17, 22, respectively for the four data sets that can be handled. These results confirm with Theorem 3, claiming that the ML solution of our estimation problem is sparse.

The average computational time for the MM method to solve the GP hyper-parameters in one MC run is around 1 minute, 25 minutes, 10 minutes, 9 minutes, respectively for the four smaller data sets that can be handled. From next subsection on, we will solely focus on the new 1-D GSM kernel with much reduced model complexity.

C. Performance of the 1-D GSM Kernel with MM Method

In the previous subsection, we showed the performance of the GSM kernel with 2-D grids. The model complexity, m , is expected to be large for good performance. We need to reduce the model complexity. We resort to the GSM kernel with 1-D grids as given in Eq.(11), for which we sample m frequency parameters μ_i , $i = 1, 2, \dots, m$ uniformly from the given frequency region $[0, 1/2)$, while fix the variance parameter to a small constant, $\sigma = 0.001$. The GP hyper-parameters are solved via the sequential MM method, for which the initial guess of α_i , is first generated from a Gaussian distribution with zero mean and large variance, say $\sigma_\alpha^2 = 10$, and then finalized by $\max(\alpha_i, 0)$. Random restart is not used for fair comparison.

To shed some light on its performance, we let $m = 100, 200, 300, 400, 500$ and repeat the tests as conducted in the previous subsection for each m . We compare the prediction MSE obtained by the 2-D GSM kernel with 20,000 grids randomly sampled from the 2-D (μ, σ) -space and the 1-D GSM kernel with only 500 grids uniformly selected from the 1-D μ -space (with a fixed $\sigma = 0.001$). The results are shown in Table III. In total 100 independent MC runs were conducted

TABLE III

PREDICTION MSE GENERATED BY TWO GSM KERNELS (ONE IS USING $m = 20000$ 2-D GRIDS VS. THE OTHER USING $m = 500$ 1-D GRIDS).

| Name | 1-D MSE | 1-D Iterations | 1-D PFR | 2-D MSE | 2-D Iterations |
|-------------|---------|----------------|---------|---------|----------------|
| ECG | 1.3E-02 | 24 | 0.01 | NA | NA |
| CO2 | 1.5E+00 | 10 | 0.17 | NA | NA |
| Electricity | 4.7E+03 | 2 | 0.07 | 6.8E+03 | 2 |
| Employment | 1.1E+02 | 23 | 0.06 | 3.9E+01 | 14 |
| Hotel | 8.9E+02 | 14 | 0.02 | 2.4E+03 | 6 |
| Passenger | 1.9E+02 | 28 | 0.02 | 1.7E+02 | 13 |
| Clay | 1.9E+02 | 25 | 0.12 | NA | NA |
| Unemploy. | 3.6E+03 | 9 | 0.10 | NA | NA |

TABLE IV

PREDICTION MSE OF THE GSMGP WITH $m = 500$ 1-D GRIDS VS. SMGP WITH $Q = 500$ GAUSSIAN MODES.

| Name | GSMGP MSE | GSMGP CT (s) | SMGP MSE | SMGP CT (s) | SMGP PFR |
|-------------|-----------|--------------|----------|-------------|----------|
| ECG | 1.3E-02 | 140.4 | 1.9E-02 | 3.4E+03 | 0.3 |
| CO2 | 1.5E+00 | 69.3 | 1.1E+00 | 2.0E+03 | 0.07 |
| Electricity | 4.7E+03 | 1.46 | 7.5E+03 | 1.0E+02 | 0 |
| Employment | 1.1E+02 | 31.2 | 0.7E+02 | 2.5E+02 | 0.03 |
| Hotel | 8.9E+02 | 17.5 | 2.8E+03 | 2.8E+02 | 0.97 |
| Passenger | 1.9E+02 | 14.7 | 1.6E+02 | 1.1E+02 | 0.23 |
| Clay | 1.9E+02 | 140.4 | 3.3E+02 | 3.4E+03 | 0 |
| Unemploy. | 3.6E+03 | 42.3 | 1.4E+04 | 1.4E+03 | 0.57 |

to compute the program fail rate as well as the prediction MSE after excluding the meaningless estimates. To better visualize the results, we show the training and prediction performance of the resulting GSMGP on the *Electricity* and *Unemployment* data sets in one specific MC run in Fig. 2. Similar results for all data sets are given in the supplement.

Some observations from our experimental results are as follows. First, in a majority of cases, the prediction MSE generated by the 1-D GSM kernel degrades slightly as compared to that generated by the 2-D GSM kernel. This result is not surprising as the latter better covers the parameter space. On the other hand, the 2-D GSM kernel may overfit the training data in some cases, as seen for the *Electricity* data set in Table III. Second, due to the significantly reduced model complexity, the MM method can handle much longer time series with the 1-D GSM kernel. Although the number of iterations required by the MM method increases for the 1-D GSM kernel, the overall computational time for the eight data sets has been reduced significantly from several minutes to several seconds. The surprisingly low computational time is also due to the low-rank property of all GSM sub-kernel matrices with $\sigma = 0.001$, supported by Theorem 1. Lastly, although not shown in the Table, the performance of the 1-D GSM kernel becomes better and more stable as m increases to around 500 grids but further increment would not help much for the selected data sets. Moreover, in Table V, we compare the GSMGP with an upgraded SMGP with $Q = 500$, which leads to much improved prediction MSE and more stable numerical solution than that of $Q = 10$ Gaussian modes, however at the cost of much longer computational time. But still for some data sets, e.g., the *Unemployment* and *Hotel*,

⁴Specifications: Intel(R) Core(TM) i7-8700 CPU 3.2GHz, 3192MHz, 6 cores, 16GB RAM with MATLAB2017a installed

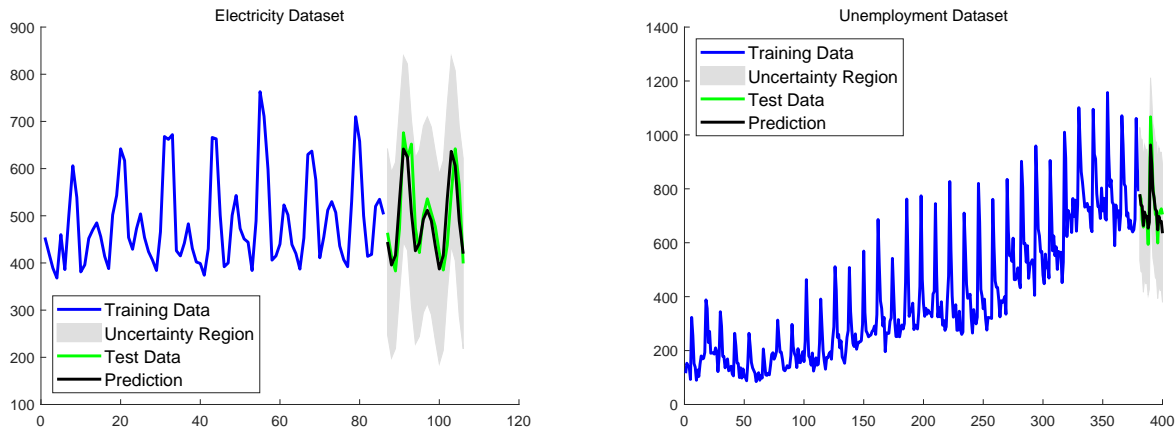


Fig. 2. Training and test performance of the GSMGP using 1-D GSM kernel with $\sigma = 0.001$ and $m = 500$ uniformly generated grids. The optimal weights are solved via the MM method.

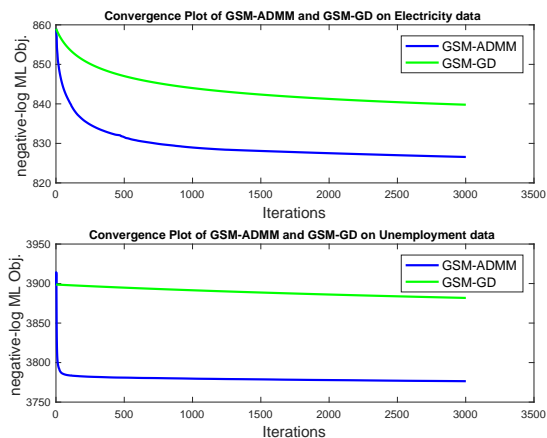


Fig. 3. Negative log-likelihood versus iterations of the proposed ADMM as compared to the classic gradient projection.

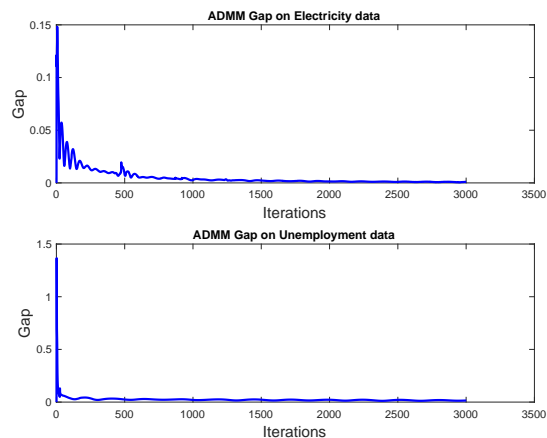


Fig. 4. Gap of the equality constraint, $\|S^k C^k - I\|_F$ versus the iterations of the proposed ADMM.

the SMGP gets stuck at bad local minimal more frequently than our GSMGP. As a summary, the new 1-D GSM kernel has achieved overall better prediction results with much less reduced computational time and higher stability as compared to the original SM kernel with a large number of Gaussian modes.

In the following, we show the benefits of using Nyström to further speed up the computations. We still stick to the GSM kernel with $m = 500$ 1-D grids and fixed $\sigma = 0.001$. We randomly sample only 5% percent of the complete training inputs for constructing a Nyström approximation of every sub-kernel matrix $\mathbf{K}_i, i = 1, 2, \dots, m$. The results in Table V show the prediction MSE as well as the computational time that the MM method requires to converge in one particular MC run initialized with all zeros. The total computation time is not reduced much in this case because the sub-kernel matrices have low-rank (refer to the fifth property of the GSM kernel as given by Eq.(11)) and this nice property matches perfectly with the MM method according to our remark 3 given in Section V. For all data sets, we computed the rank of all

sub-kernel matrices numerically and recorded the maximum rank, the minimum rank and the mean rank in Table VI which demonstrate $\max_i \text{rank}\{\mathbf{K}_i\} \approx 2 \min_i \text{rank}\{\mathbf{K}_i\} \approx 1.3\sqrt{n}$; the mean rank is fairly close to the maximum rank because most of the sub-kernels have rank close to the maximum rank. Therefore, in light of the remark 3 of Section V, the computational complexity of the MM method is approximately $\mathcal{O}(mn^{3/2})$ instead of the worst case $\mathcal{O}(mn^3)$. When we handle longer time series, random Fourier feature approximation may help save more memory while maintain similar RAE, as was shown in the supplement.

In all above experiments, we use the default setup of the 1-D GSM kernel, which is very simple to use. But as we pointed out in Section III, using the nonparametric Welch periodogram of the data to guide an advanced setup may be beneficial in various aspects. Due to space limitation, we show the periodogram of each data set versus the spectral density constructed using the optimal weights obtained for one specific MC in the supplement. As we can see, the periodogram indeed provides rich information for configuring the GSM kernel and

TABLE V
PREDICTION MSE GENERATED BY THE 1-D GSM KERNEL VERSUS ITS NYSTRÖM APPROXIMATION, SHORT AS NY-GSM.

| Name | GSM MSE | GSM CT | NY-GSM MSE | NY-GSM CT |
|-------------|---------|--------|------------|-----------|
| ECG | 1.3E-02 | 122s | 1.3E-02 | 116s |
| CO2 | 9.3E-01 | 24s | 9.3E-01 | 22s |
| Electricity | 3.0E+03 | 0.9s | 3.0E+03 | 0.2s |
| Employment | 6.8E+01 | 12s | 6.8E+01 | 5s |
| Hotel | 4.3E+02 | 3s | 4.3E+02 | 1s |
| Passenger | 2.4E+02 | 8s | 2.9E+02 | 3s |
| Clay | 8.5E+01 | 60s | 8.5E+01 | 50s |
| Unemploy. | 2.3E+03 | 8s | 2.3E+03 | 3s |

TABLE VI
MAXIMUM RANK, MINIMUM RANK, AND MEAN RANK OF THE SELECTED $m = 500$ GSM SUB-KERNEL MATRICES USED IN THE ABOVE EXPERIMENTS.

| Name | max rank GSM sub-kernels | min rank sub-kernels | mean rank sub-kernels |
|--------------|-----------------------------|-------------------------|--------------------------|
| ECG | 34 | 17 | 33 |
| CO2 | 27 | 13 | 25 |
| Electricity | 14 | 7 | 13 |
| Employment | 16 | 8 | 15 |
| Hotel | 14 | 7 | 13 |
| Passenger | 14 | 7 | 13 |
| Clay | 26 | 13 | 25 |
| Unemployment | 24 | 12 | 23 |

optimizing its associated hyper-parameters.

D. Performance of the 1-D GSM Kernel with ADMM

In section V-B, we introduced a nonlinearly constrained ADMM, short as GSM-ADMM, for optimizing the hyper-parameters of the 1-D GSM kernel, i.e., the weights α . In the following experiments, we aim to compare it with other two numerical methods, namely the classic gradient projection (details see our supplement) and the sequential MM method, short as GSM-GD and GSM-MM, respectively. The performance is measured in terms of the objective function value and the prediction MSE.

We conduct some experiments on a small data set and a moderate data set, as the proposed method is not suitable for big data set due to the $\mathcal{O}(n^3)$ complexity. To keep alignment with the previous experiments, we stick to the 1-D GSM kernel with $m = 500$ grids uniformly sampled from $[0, 1/2)$ and fixed $\sigma = 0.001$.

The algorithmic setup of our nonlinearly constrained ADMM as given in Algorithm 2 are in order. To update S , we let $It_S = 1000$, $\epsilon_S = 10^{-15}$, $\delta = 1$. For selecting the step size in light of the Armijo rule, we let $s = 10^{-4}$, $\beta = 1/5$, $h = 10^{-5}$. The remainders are $\rho = 100$, $\rho' = \rho/2 = 50$, $\epsilon_{ADMM} = 10^{-3}$.

As for the initial guess, we let $\Lambda^{(0)} = I$, for the *Electricity* data set; $\alpha^{(0)}$ is obtained by fitting the nonparametric Welch periodogram via the L_1 -norm regularized least-squares mentioned in Section III, while for the *Unemployment* data set, $\alpha^{(0)}$ is obtained by running just one iteration of the sequential MM method. The same initial guesses were applied to the GSM-GD for fair comparisons. The experimental results are

summarized in Table VII. Instead of striving to find a local

TABLE VII
PERFORMANCE OF THREE NUMERICAL OPTIMIZATION METHODS IN TERMS OF THE OBJECTIVE FUNCTION VALUE, THE PREDICTION MSE, AND THE COMPUTATIONAL TIME

| Performance Metric | Electricity | Unemployment |
|--------------------|-------------|--------------|
| GSM-GD Objective | 8.330E+02 | 3.838E+03 |
| GSM-MM Objective | 8.284E+02 | 3.779E+03 |
| GSM-ADMM Objective | 8.266E+02 | 3.776E+03 |
| GSM-GD MSE | 4.426E+03 | 1.481E+04 |
| GSM-MM MSE | 3.037E+03 | 2.248E+03 |
| GSM-ADMM MSE | 2.220E+03 | 2.222E+03 |
| GSM-GD CT (s) | 2272s | 79189s |
| GSM-MM CT (s) | 0.93s | 8.40s |
| GSM-ADMM CT (s) | 6351.17s | 160367.25s |

minimum, we restrict the maximum number of iterations of the nonlinearly constrained ADMM due to its relatively slow convergence rate. Although the ADMM has not converged yet, it already found a weight estimate α that leads to the smallest objective function value and prediction MSE among the selected numerical methods. However, the proposed nonlinearly constrained ADMM is less favorable than the MM method in terms of the computational time in both cases.

As yet another comparison between the GSM-GD and GSM-ADMM, we showed the negative log-likelihood value versus iterations in Fig. 3. It is clear from the results that the GSM-ADMM shows faster convergence rate as compared to GSM-GD. The gap of new introduced equality constraint is also depicted versus iterations in Fig. 4.

Lastly, we give some guidance on the selection of a few key parameters of the proposed ADMM:

- Regularization parameter ρ : In general, a smaller ρ leads to faster convergence rate of the method, while a larger ρ leads to more stable convergence progress and smaller gap in the equality constraint but at a slower convergence rate. Good trade-off needs to be taken care of.
- Tolerance ϵ_S : in our ADMM, ϵ_S is typically chosen small to waive the effect of the inexact solution of the sub-problem in Eq.(28) and improve the overall convergence performance. However, a too small ϵ_S , on the other hand, is prohibited due to the high computational cost required for function evaluations. As rule-of-thumb, we could choose $\epsilon_S \in [10^{-10}, 10^{-15}]$.

VII. CONCLUSION AND OUTLOOK

We studied automatic, optimal stationary kernel design with the good aim to let data choose the most appropriate kernel. We modified the SM kernel in the frequency domain by fixing the frequency and variance parameters to a big number of pre-selected grids. We conducted thorough studies on the properties of the resultant 1-D GSM kernel, including the sampling strategies of the grids, validity and low-rank property of all sub-kernels, and user-friendly initialization. The resultant GSM kernel demonstrates itself to be a linear multiple kernel. The ML based hyper-parameter optimization problem falls in difference-of-convex program and the solution is widely known to be sparse. Experimental results showed

that the MM method achieved the best overall performance in various aspects, including convergence speed, economical computational time, insensitivity to an initial guess, competent fitting and prediction performance, etc. The fast computational speed of the MM method is obtained due to the low-rank properties of all GSM sub-kernels. On the other hand, the proposed ADMM showed great potential to achieve better local minimum but at the cost of larger computational time. Experimental results based on various classic time series data sets confirmed that the proposed 1-D GSM kernel is able to generate overall better performance than its 2-D counterpart and several other salient competitors of similar kind. Although the proposed 1-D GSM kernel showed outstanding prediction performance and very fast computational speed, it is more favorable to be used for low-dimensional time series.

ACKNOWLEDGEMENT

The author Feng Yin would like to thank Prof. Abdelhak M. Zoubir from Technische Universität Darmstadt and Prof. Xiaodong LI from UC Davis for the fruitful discussions on this manuscript during their visits at CUHK(SZ).

APPENDIX

A. Proof of property (1): GSM kernel is a valid kernel

A necessary and sufficient condition for a function $k(\mathbf{x}, \mathbf{x}')$ to be a valid kernel according to [42] is that the corresponding kernel matrix, whose (i, j) -th entry is given by $k(\mathbf{x}_i, \mathbf{x}_j)$, is PSD for all possible choices of $\mathbf{x} \in \mathcal{X}$.

For the proof, we need the following fundamental operations for constructing a new valid kernel $k(\mathbf{x}, \mathbf{x}')$ that are well known from [42] and [24]:

$$k(\mathbf{x}, \mathbf{x}') = f(\mathbf{x})f(\mathbf{x}') \quad (36a)$$

$$k(\mathbf{x}, \mathbf{x}') = ck_1(\mathbf{x}, \mathbf{x}') \quad (36b)$$

$$k(\mathbf{x}, \mathbf{x}') = f(\mathbf{x})k_1(\mathbf{x}, \mathbf{x}')f(\mathbf{x}') \quad (36c)$$

$$k(\mathbf{x}, \mathbf{x}') = \exp(k_1(\mathbf{x}, \mathbf{x}')) \quad (36d)$$

$$k(\mathbf{x}, \mathbf{x}') = k_1(\mathbf{x}, \mathbf{x}') + k_2(\mathbf{x}, \mathbf{x}') \quad (36e)$$

$$k(\mathbf{x}, \mathbf{x}') = k_1(\mathbf{x}, \mathbf{x}') \cdot k_2(\mathbf{x}, \mathbf{x}') \quad (36f)$$

where $k_1(\mathbf{x}, \mathbf{x}')$ and $k_2(\mathbf{x}, \mathbf{x}')$ are both known valid kernels; $f(\mathbf{x}) : \mathbb{R}^d \rightarrow \mathbb{R}$ is any function; $c \geq 0$ is a constant. In our work, $\mathbf{x} = t$ and $d = 1$.

We will use the above results to prove that each sub-kernel function (omitting the subscript i) $k(t, t'; \sigma^2, \mu) = \exp[-2\pi^2(t-t')^2\sigma^2] \cos(2\pi(t-t')\mu)$ is a valid kernel. First, we let $k(t, t'; \sigma^2, \mu) = k_1(t, t'; \sigma^2) \cdot k_2(t, t'; \mu)$, where $k_1(t, t'; \sigma^2) \triangleq \exp[-2\pi^2(t-t')^2\sigma^2]$ and $k_2(t, t'; \mu) \triangleq \cos(2\pi(t-t')\mu)$. The first part $k_1(t, t'; \sigma^2)$ can be reformulated as

$$k_1(t, t'; \sigma^2) = \exp[-2\pi^2\sigma^2 t^2] \exp[4\pi^2\sigma^2 t t'] \exp[-2\pi^2\sigma^2 t'^2] \\ = f_1(t) \exp[4\pi^2\sigma^2 \cdot k_{11}(t, t')] f_1(t'), \quad (37)$$

where $f_1(t) \triangleq \exp[-2\pi^2\sigma^2 t^2]$ and $k_{11}(t, t') \triangleq t t'$ is the well known, valid linear kernel. Applying the fundamental operations given in Eq.(36b), Eq.(36d), and Eq.(36c) in turn, yields a valid kernel $k_1(t, t'; \sigma^2)$.

Next, we prove $k_2(t, t'; \mu) \triangleq \cos(2\pi(t-t')\mu)$ is also a valid kernel. This is done by reformulating the kernel as:

$$k_2(t, t'; \mu) = \cos(2\pi\mu t - 2\pi\mu t') \\ = \cos(2\pi\mu t) \cos(2\pi\mu t') + \sin(2\pi\mu t) \sin(2\pi\mu t') \\ = f_{21}(t)f_{21}(t') + f_{22}(t)f_{22}(t'), \quad (38)$$

where $f_{21}(t) \triangleq \cos(2\pi\mu t)$ and $f_{22}(t) \triangleq \sin(2\pi\mu t)$. Applying the fundamental operations given in Eq.(36a) and Eq.(36e) in turn, yields a valid kernel $k_2(t, t'; \mu)$.

Since both $k_1(t, t'; \sigma^2)$ and $k_2(t, t'; \mu)$ are valid kernels, according to Eq.(36f), $k(t, t'; \sigma^2, \mu)$ is a valid kernel. The above proof holds generally for any $\sigma^2 \in \mathbb{R}_+$ and $\mu \in \mathbb{R}$. Consequently, each sub-kernel matrix \mathbf{K}_i is a PSD matrix and so is $\mathbf{K} = \sum_{i=1}^m \alpha_i \mathbf{K}_i$.

B. Verification of property (5): low-rank property

We let \mathbf{K}_i for a given grid with (μ_i, σ^2) be the $n \times n$ kernel matrix of the i -th sub-kernel $k_i(t, t')$ given in Eq.(11) with t and t' in $\{1, 2, \dots, n\}$. The kernel matrix can be expressed in the form of Hadamard product as $\mathbf{K}_i = \mathbf{K}_i^{exp} \circ \mathbf{K}_i^{cos}$. Here, \mathbf{K}_i^{exp} can be seen as the kernel matrix of its corresponding stationary kernel function $k^{exp}(\tau) \triangleq \exp(-2\pi^2\tau^2\sigma^2)$ and \mathbf{K}_i^{cos} can be seen as the kernel matrix of its corresponding stationary kernel function $k_i^{cos}(\tau) \triangleq \cos(2\pi\tau\mu_i)$. According to the rank inequality of Hadamard product of two matrices [43], we have

$$rank(\mathbf{K}_i) \leq rank(\mathbf{K}_i^{exp}) \cdot rank(\mathbf{K}_i^{cos}). \quad (39)$$

We need the following two lemmas, (1) $rank(\mathbf{K}_i^{cos}) = 2$ for any grid i and (2) $rank(\mathbf{K}_i^{exp}) \ll \frac{n}{2}$ for sufficiently small σ^2 .

Lemma 1. For the kernel function $k_i^{cos}(\tau) = \cos(2\pi\tau\mu_i)$ with any $\mu_i \in (0, 1/2)$, the rank of the corresponding kernel matrix is always equal to 2, i.e., $rank(\mathbf{K}_i^{cos}) = 2$ for any i .

Proof. The proof is as follows. It is obvious that first column of \mathbf{K}_i^{cos} , denoted by $\mathbf{k}_1 = [\cos(0x), \cos(1x), \dots, \cos((n-1)x)]^T$ and the second column, denoted by $\mathbf{k}_2 = [\cos(-1x), \cos(0x), \cos(1x), \dots, \cos((n-2)x)]^T$, where $x = 2\pi\mu_i$ is a constant in $(0, \pi)$ for any given μ_i , are linearly independent. While from the 3rd column onward, each column can be expressed as a linear combination of the previous two columns simply because it holds for any $j \in \{-(n-1) : 1 : (n-1)\}$, $\cos(jx) = \alpha \cos((j+2)x) + \beta \cos((j+1)x)$, where $\alpha = -1$ and $\beta = \sin(2x)/\sin(x)$ are both irrespective of j . The derivation of α and β is due to

$$\cos(jx) = (\alpha \cos(2x) + \beta \cos(x)) \cos(jx) \\ - (\alpha \sin(2x) + \beta \sin(x)) \sin(jx). \quad (40)$$

Then, we let $\alpha \cos(2x) + \beta \cos(x) = 1$ and $\alpha \sin(2x) + \beta \sin(x) = 0$, then solve for α and β . The above steps prove that the kernel matrix \mathbf{K}_i^{cos} is always of rank 2. \square

Lemma 2. For a time series with n samples, i.e., $t = 1, 2, \dots, n$, when the variance parameter is selected to be $\sigma^2 \leq \frac{2r+1}{2\pi^2(n-1)^2 \cdot C} \ll \frac{(n/2)+1}{2\pi^2(n-1)^2 \cdot C} \approx \frac{1}{4\pi^2(n-1) \cdot C}$, the rank of the kernel matrix \mathbf{K}^{exp} corresponding to $k^{exp}(\tau) = \exp(-2\pi^2\tau^2\sigma^2)$, for any $\tau \in \{0, 1, 2, \dots, n-1\}$, satisfies

$\text{rank}(\mathbf{K}^{exp}) \leq (2r + 1) \ll n/2$ for some large constant number C .

Proof. First of all, we show that there exists certain K such that for each $\tau \in \{0, 1, 2, \dots, n-1\}$, the exponential function $\exp(-2\pi^2\sigma^2\tau^2)$, short for $\exp(a\tau^2)$, can be approximated by the first K terms of its Taylor expansion, namely, $\exp(a\tau^2) = 1 + a\tau^2 + \frac{(a\tau^2)^2}{2!} + \frac{(a\tau^2)^3}{3!} + \dots + \frac{(a\tau^2)^K}{K!} + R_{K+1}$, where the remainder $R_{K+1} = \frac{(a\tau^2)^{(K+1)}}{(K+1)!} \exp(t \cdot a\tau^2)$ with $0 < t < 1$. It is known that $\lim_{K \rightarrow \infty} |R_{K+1}| \leq \exp(|a\tau^2|) \cdot \lim_{K \rightarrow \infty} \left| \frac{(a\tau^2)^{(K+1)}}{(K+1)!} \right| = 0$, hence for any $\epsilon > 0$, we can find a K such that the approximation error $|R_{K+1}| < \epsilon$. In order to give a practical guidance on the selection of σ , we aim to find a number K such that $\frac{(2\pi^2\sigma^2\tau^2)^K}{K!} > C \cdot \frac{(2\pi^2\sigma^2\tau^2)^{K+1}}{(K+1)!}$, $\forall \tau \in \{0, 1, 2, \dots, n-1\}$, where C is a large constant number. Conservatively for $\tau = n-1$, we have $\sigma^2 < \frac{K+1}{2\pi^2(n-1)^2 \cdot C}$, implying that for a fixed n , when σ shrinks, the above inequality could still hold with smaller K . Since a drastic decrease in the absolute values of consecutive terms is our indicator for good approximation using Taylor expansion, in order to achieve $K = 2r \ll n/2$, we need to select $\sigma^2 \leq \frac{2r+1}{2\pi^2(n-1)^2 \cdot C} \ll \frac{(n/2)+1}{2\pi^2(n-1)^2 \cdot C} \approx \frac{1}{4\pi^2(n-1) \cdot C}$.

Next, we show that the rank of \mathbf{K}^{exp} is at most $(2r+1)$ for a sufficiently small σ^2 due to the fact that $\exp(a\tau^2)$ can be well approximated by a linear combination of its $(2r+1)$ previous terms $\exp(a(\tau+1)^2), \exp(a(\tau+2)^2), \dots, \exp(a(\tau+2r+1)^2)$ regardless of τ . Our proof is as follows. For any give $\tau \in \{1, 2, \dots, n-1\}$, we have the approximation

$$\begin{aligned} \exp(a\tau^2) &\approx 1 + a\tau^2 + \frac{(a\tau^2)^2}{2!} + \frac{(a\tau^2)^3}{3!} + \dots + \frac{(a\tau^2)^r}{r!} \\ &= \tilde{a}_{0,1} + \tilde{a}_{0,1}\tau + \tilde{a}_{0,2}\tau^2 + \dots + \tilde{a}_{0,2r}\tau^{2r}, \end{aligned} \quad (41)$$

where some of the coefficients are zeros. Similarly, for the i -th previous term we have $\exp(a(\tau+i)^2) \approx \tilde{a}_{i,0} + \tilde{a}_{i,1}\tau + \tilde{a}_{i,2}\tau^2 + \dots + \tilde{a}_{i,2r}\tau^{2r}$ for any $i \in \{1, 2, \dots, 2r+1\}$. With the introduction of a coefficient matrix $\tilde{\mathbf{A}}_{(2r+1) \times (2r+1)}$ whose ij -th element is $\tilde{a}_{i,j-1}$, $\tilde{\mathbf{a}}_0 = [1, \tilde{a}_{0,1}, \tilde{a}_{0,2}, \dots, \tilde{a}_{0,2r}]^T$ and $\boldsymbol{\beta} \in \mathbb{R}^{(2r+1)}$, we can construct a linear system $\tilde{\mathbf{A}}\boldsymbol{\beta} = \tilde{\mathbf{a}}_0$. Solving this linear system yields $\exp(a\tau^2) = [\exp(a(\tau+1)^2), \exp(a(\tau+2)^2), \dots, \exp(a(\tau+2r+1)^2)]\boldsymbol{\beta}$. An important fact is that the solution of $\boldsymbol{\beta}$ has nothing to do with τ , since both $\tilde{\mathbf{A}}$ and $\tilde{\mathbf{a}}_0$ are only in terms of the fixed $2\pi\sigma^2$. As a result, for any $j > 2r+1$ and $2r \ll n/2$, the j -th column of \mathbf{K}^{exp} can be reproduced by a linear combination of its $(2r+1)$ previous columns.

Combining the above two parts completes the proof of this theorem. \square

Corollary 2. *Following the above lemma, when $\sigma \rightarrow 0$, $\text{rank}(\mathbf{K}^{exp}) \rightarrow 1$.*

C. Derivation of Eq.(34)

Solving Eq.(33) for updated α_i yields

$$\begin{aligned} \alpha_i^{k+1} &= \frac{-tr \left[\left(\tilde{\mathbf{K}}_{-i} \mathbf{K}_i + \sigma_e^2 \mathbf{K}_i \right) \mathbf{S}^{k+1,T} \mathbf{S}^{k+1} - \mathbf{S}^{k+1} \mathbf{K}_i \right]}{tr \left(\mathbf{K}_i^T \mathbf{S}^{k+1,T} \mathbf{S}^{k+1} \mathbf{K}_i \right)} \\ &\quad - \frac{tr \left(\boldsymbol{\Lambda}^{k,T} \mathbf{S}^{k+1} \mathbf{K}_i \right)}{\rho \cdot tr \left(\mathbf{K}_i^T \mathbf{S}^{k+1,T} \mathbf{S}^{k+1} \mathbf{K}_i \right)}. \end{aligned} \quad (42)$$

It is noted that $\mathbf{K}_i^T = \mathbf{K}_i$ due to the symmetric property. Replace $\left(\tilde{\mathbf{K}}_{-i} \mathbf{K}_i + \sigma_e^2 \mathbf{K}_i \right)$ with $\left(\tilde{\mathbf{C}}_i - \alpha_i^k \mathbf{K}_i \right) \mathbf{K}_i^T$ in the above equation gives

$$\begin{aligned} \alpha_i^{k+1} &= \frac{-tr \left[\mathbf{K}_i^T \mathbf{S}^{k+1,T} \mathbf{S}^{k+1} \left(\tilde{\mathbf{C}}_i - \alpha_i^k \mathbf{K}_i \right) \right]}{tr \left(\mathbf{K}_i^T \mathbf{S}^{k+1,T} \mathbf{S}^{k+1} \mathbf{K}_i \right)} \\ &\quad + \frac{tr \left(\mathbf{K}_i^T \mathbf{S}^{k+1,T} \left(\mathbf{I}_n - \frac{1}{\rho} \boldsymbol{\Lambda}^k \right) \right)}{tr \left(\mathbf{K}_i^T \mathbf{S}^{k+1,T} \mathbf{S}^{k+1} \mathbf{K}_i \right)}. \end{aligned} \quad (43)$$

Dragging $-\alpha_i \mathbf{K}_i$ outside of the first term, merging the other terms, and using the fact that both \mathbf{K}_i and \mathbf{S} are symmetric, yields Eq. (34). When $\mathbf{S}^{k+1} \tilde{\mathbf{C}}_i$ is close to \mathbf{I}_n , the above update is approximately

$$\alpha_i^{k+1} \approx \alpha_i^k - \frac{1}{\rho} \cdot \frac{tr \left[\mathbf{K}_i \mathbf{S}^{k+1} \boldsymbol{\Lambda}^k \right]}{tr \left(\mathbf{K}_i \mathbf{S}^{k+1} \mathbf{S}^{k+1} \mathbf{K}_i \right)}. \quad (44)$$

REFERENCES

- [1] F. Yin, X. He, L. Pan, T. Chen, Z.-Q. Luo, and S. Theodoridis, "Sparse structure enabled grid spectral mixture kernel for temporal Gaussian process regression," in *Proceedings of International Conference on Information Fusion*, Cambridge, UK, July 2018, pp. 47–54.
- [2] C. E. Rasmussen and C. I. K. Williams, *Gaussian Processes for Machine Learning*. MIT Press, 2006.
- [3] J. Lee, J. Sohl-dickstein, J. Pennington, R. Novak, S. Schoenholz, and Y. Bahri, "Deep neural networks as Gaussian processes," in *Proceedings of International Conference on Learning Representations*, Vancouver, BC, Canada, 2018.
- [4] A. G. de G. Matthews, J. Hron, M. Rowland, R. E. Turner, and Z. Ghahramani, "Gaussian process behaviour in wide deep neural networks," in *Proceedings of International Conference on Learning Representations*, Vancouver, BC, Canada, 2018.
- [5] J. Snoek, H. Larochelle, and R. P. Adams, "Practical Bayesian optimization of machine learning algorithms," in *Proceedings of International Conference on Neural Information Processing Systems*, Lake Tahoe, Nevada, US, 2012, pp. 2951–2959.
- [6] S. Theodoridis, *Machine Learning: A Bayesian and Optimization Perspective*. Academic Press, 2015.
- [7] Y. Zhao, F. Yin, F. Gunnarsson, F. Hultkratz, and J. Fagerlind, "Gaussian processes for flow modeling and prediction of positioned trajectories evaluated with sports data," in *Proceedings of International Conference on Information Fusion*, Heidelberg, Germany, July 2016, pp. 1461–1468.
- [8] J. Han, X. Zhang, and F. Wang, "Gaussian process regression stochastic volatility model for financial time series," *IEEE Journal of Selected Topics in Signal Processing*, vol. 10, no. 6, pp. 1015–1028, September 2016.
- [9] A. Prakash, S. Xu, R. Rajagopal, and H.-Y. Noh, "Robust building energy load forecasting using physically-based kernel models," *Energies*, vol. 11, no. 4, pp. 1–21, 2018.
- [10] M. Gönen and E. Alpaydin, "Multiple kernel learning algorithms," *Journal of Machine Learning Research*, vol. 12, pp. 2211–2268, February 2011.
- [11] G. Lanckriet, N. Cristianini, P. L. Bartlett, L. E. Ghaoui, and M. I. Jordan, "Learning the kernel matrix with semi-definite programming," *Journal of Machine Learning Research*, vol. 5, pp. 27–72, December 2004.
- [12] A. Aravkin, J. V. Burke, A. Chiuso, and G. Pillonetto, "Convex vs non-convex estimators for regression and sparse estimation: the mean squared error properties of ARD and GLasso," *Journal of Machine Learning Research*, vol. 15, pp. 217–252, 2014.

- [13] J. Zhuang, J. Wang, C. H. Hoi, and X. Lan, "Unsupervised multiple kernel learning," in *Proceedings of the Asian Conference on Machine Learning*, vol. 20, Taiwan, November 2011, pp. 129–144.
- [14] R. Senanayake, S. O'Callaghan, and F. Ramos, "Predicting spatio-temporal propagation of seasonal influenza using variational Gaussian process regression," in *Proceedings of AAAI Conference on Artificial Intelligence*, Phoenix, Arizona, US, 2016, pp. 3901–3907.
- [15] Y. Xu, W. Xu, F. Yin, J. Lin, and S. Cui, "High-accuracy wireless traffic prediction: A GP-based machine learning approach," in *Proceedings of IEEE Global Communications Conference*, Singapore, 2017, pp. 1–6.
- [16] M. Lázaro-Gredilla, J. Quiñero Candela, C. E. Rasmussen, and A. R. Figueiras-Vidal, "Sparse spectrum Gaussian process regression," *Journal of Machine Learning Research*, vol. 11, pp. 1865–1881, August 2010.
- [17] D. Duvenaud, J. R. Lloyd, R. Grosse, J. B. Tenenbaum, and Z. Ghahramani, "Structure discovery in nonparametric regression through compositional kernel search," in *Proceedings of International Conference on Machine Learning*, Atlanta, USA, 2013, pp. 1166–1174.
- [18] A. Wilson and R. P. Adams, "Gaussian process kernels for pattern discovery and extrapolation," in *Proceedings of International Conference on Machine Learning*, Atlanta, USA, 2013, pp. 1067–1075.
- [19] A. G. Wilson, "Covariance kernels for fast automatic pattern discovery and extrapolation with Gaussian processes," Ph.D. dissertation, University of Cambridge, UK, 2014.
- [20] K. Krauth, E. Bonilla, C. K., and F. M., "AutoGP: Exploring the capabilities and limitations of Gaussian process models," in *Proceedings of Conference on Uncertainty in Artificial Intelligence*, Sydney, Australia, August 2017.
- [21] R. M. Neal, "Monte Carlo implementation of Gaussian process models for Bayesian regression and classification," Dept. of Statistics, University of Toronto, Canada, Tech. Rep. Technical Report No. 9702, 1997.
- [22] C. Wang and R. M. Neal, "MCMC methods for Gaussian process models using fast approximations for the likelihood," University of Toronto, Canada, Tech. Rep., 2013.
- [23] N. I. Achieser, *Theory of Approximation*. Dover Publications, Inc., New York, 1992.
- [24] C. Bishop, *Machine Learning and Pattern Recognition*. Springer, 2006.
- [25] J. R. Gilbert, C. Moler, and R. Schreiber, "Sparse matrices in MATLAB: Design and implementation," *SIAM J. Matrix Anal. Appl.*, vol. 13, no. 1, pp. 333–356, January 1992.
- [26] B. W. Silverman, *Density Estimation for Statistical and Data Analysis*. London: Chapman and Hall, 1986.
- [27] A. V. Oppenheim and R. W. Schaffer, *Digital Signal Processing*. Englewood Cliffs, N.J.: Prentice-Hall, 1993.
- [28] S.-J. Kim, K. Koh, M. Lustig, S. Boyd, and D. Gorinevsky, "A method for large-scale L1-regularized least squares," *IEEE Journal on Selected Topics in Signal Processing*, vol. 1, pp. 606–617, 2007.
- [29] C. Williams and M. Seeger, "Using the Nyström method to speed up kernel machines," in *Advances in Neural Information Processing Systems*. MIT Press, 2001, pp. 682–688.
- [30] A. Rahimi and B. Recht, "Random features for large-scale kernel machines," in *Proceedings of International Conference on Neural Information Processing Systems*, Vancouver, British Columbia, Canada, 2007, pp. 1177–1184.
- [31] P. Bouboulis, S. Chouvardas, and S. Theodoridis, "Online distributed learning over networks in RKH spaces using random Fourier features," *IEEE Transactions on Signal Processing*, vol. 66, pp. 1920–1932, April 2018.
- [32] Q. Le, T. Sarlos, and A. J. Smola, "Fastfood—computing hilbert space expansions in loglinear time," in *Proceedings of International Conference on Machine Learning*, Atlanta, USA, 2013, pp. 244–252.
- [33] S. Boyd and L. Vandenberghe, *Convex Optimization*. Cambridge University Press, 2004.
- [34] T. Chen, M. S. Anderson, and L. Ljung, "System identification via sparse multiple kernel-based regularization using sequential convex optimization techniques," *IEEE Transactions on Automatic Control*, vol. 59, no. 11, pp. 2933–2945, 2014.
- [35] T. Lipp and S. Boyd, "Variations and extensions of the convex-concave procedure," *Optimization and Engineering*, vol. 17, no. 2, pp. 263–287, 2016.
- [36] S. Boyd, N. Parikh, E. Chu, B. Peleato, and J. Eckstein, "Distributed optimization and statistical learning via the alternating direction method of multipliers," *Foundation Trends Machine Learning*, vol. 3, no. 1, pp. 1–122, January 2011.
- [37] M. Grant and S. Boyd, "CVX: Matlab software for disciplined convex programming, version 2.1," <http://cvxr.com/cvx>, 2014.
- [38] M. Hong, Z. Q. Luo, and M. Razaviyayn, "Convergence analysis of alternating direction method of multipliers for a family of nonconvex problems," *SIAM Journal on Optimization*, vol. 26, no. 1, pp. 337–364, January 2016.
- [39] D. P. Bertsekas, *Nonlinear Programming, 3rd. Edition*. Athena Scientific, Belmont, Mass. US., 2016.
- [40] D. P. Wipf and B. D. Rao, "Sparse Bayesian learning for basis selection," *IEEE Transactions on Signal Processing*, vol. 52, no. 8, pp. 2153–2164, August 2004.
- [41] D. Garcia, "Robust smoothing of gridded data in one and higher dimensions with missing values," *Computational Statistics and Data Analysis*, vol. 54, pp. 1167–1178, September 2010.
- [42] J. Shawe-Taylor and N. Cristianini, *Kernel Methods for Pattern Analysis*. Cambridge University Press, 2004.
- [43] G. P. H. Styan, "Hadamard products and multivariate statistical analysis," *Linear Algebra and its Applications*, vol. 6, pp. 217–240, 1973.