# Hilbert Space Methods for Reduced-Rank Gaussian Process Regression

Arno Solin · Simo Särkkä

**Abstract** This paper proposes a novel scheme for reduced-rank Gaussian process regression. The method is based on an approximate series expansion of the covariance function in terms of an eigenfunction expansion of the Laplace operator in a compact subset of $\mathbb{R}^d$. On this approximate eigenbasis the eigenvalues of the covariance function can be expressed as simple functions of the spectral density of the Gaussian process, which allows the GP inference to be solved under a computational cost scaling as $\mathcal{O}(nm^2)$ (initial) and $\mathcal{O}(m^3)$ (hyperparameter learning) with $m$ basis functions and $n$ data points. Furthermore, the basis functions are independent of the parameters of the covariance function, which allows for very fast hyperparameter learning. The approach also allows for rigorous error analysis with Hilbert space theory, and we show that the approximation becomes exact when the size of the compact subset and the number of eigenfunctions go to infinity. We also show that the convergence rate of the truncation error is independent of the input dimensionality provided that the differentiability order of the covariance function is increases appropriately, and for the squared exponential covariance function it is always bounded by $\sim 1/m$ regardless of the input dimensionality. The expansion generalizes to Hilbert spaces with an inner product which is defined as an integral over a specified input density. The method is compared to previously

A. Solin
Department of Computer Science
Aalto University
P.O. Box 15400, FI-00076 Aalto, Finland
Tel.: +358-40-5776226
E-mail: arno.solin@aalto.fi

S. Särkkä
Department of Electrical Engineering and Automation
Aalto University

proposed methods theoretically and through empirical tests with simulated and real data.

## 1 Introduction

Gaussian processes (GPs, Rasmussen and Williams, 2006) are powerful tools for non-parametric Bayesian inference and learning. In GP regression the model functions $f(\mathbf{x})$ are assumed to be realizations from a Gaussian random process prior with a given covariance function $k(\mathbf{x}, \mathbf{x}')$, and learning amounts to solving the posterior process given a set of noisy measurements $y_1, y_2, \ldots, y_n$ at some given test inputs. This model is often written in the form

$$
\begin{aligned}
f &\sim \mathcal{GP}(0, k(\mathbf{x}, \mathbf{x}')), \\
y_i &= f(\mathbf{x}_i) + \varepsilon_i,
\end{aligned}
\tag{1}
$$

where $\varepsilon_i \sim \mathcal{N}(0, \sigma_\mathrm{n}^2)$, for $i = 1, 2, \ldots, n$. One of the main limitations of GPs in machine learning is the computational and memory requirements that scale as $\mathcal{O}(n^3)$ and $\mathcal{O}(n^2)$ in a direct implementation. This limits the applicability of GPs when the number of training samples $n$ grows large. The computational requirements arise because in solving the GP regression problem we need to invert the $n \times n$ Gram matrix $\mathbf{K} + \sigma_\mathrm{n}^2 \mathbf{I}$, where $\mathbf{K}_{ij} = k(\mathbf{x}_i, \mathbf{x}_j)$, which is an $\mathcal{O}(n^3)$ operation in general.

To overcome this problem, over the years, several schemes have been proposed. They typically reduce the storage requirements to $\mathcal{O}(nm)$ and complexity to $\mathcal{O}(nm^2)$, where $m < n$. Some early methods have been reviewed in Rasmussen and Williams (2006), and Quiñonero-Candela and Rasmussen (2005b) provide a

unifying view on several methods. From a spectral point of view, several of these methods (*e.g.*, SOR, DTC, VAR, FIC) can be interpreted as modifications to the so-called *Nyström method* (see Baker, 1977; Williams and Seeger, 2001), a scheme for approximating the eigenspectrum.

For stationary covariance functions the spectral density of the covariance function can be employed: in this context the spectral approach has mainly been considered in regular grids, as this allows for the use of FFT-based methods for fast solutions (see Paciorek, 2007; Fritz et al, 2009), and more recently in terms of converting GPs to state space models (Särkkä and Hartikainen, 2012; Särkkä et al, 2013). Recently, Lázaro-Gredilla et al (2010) proposed a sparse spectrum method where a randomly chosen set of spectral points span a trigonometric basis for the problem.

The methods proposed in this article fall into the class of methods called reduced-rank approximations (see, *e.g.*, Rasmussen and Williams, 2006, Ch. 8) which are based on approximating the Gram matrix $\mathbf{K}$ with a matrix $\tilde{\mathbf{K}}$ with a smaller rank $m < n$. This allows for the use of matrix inversion lemma (Woodbury formula) to speed up the computations. It is well-known that the optimal reduced-rank approximation of the Gram (covariance) matrix $\mathbf{K}$ with respect to the Frobenius norm is $\tilde{\mathbf{K}} = \mathbf{\Phi}\mathbf{\Lambda}\mathbf{\Phi}^\mathsf{T}$, where $\mathbf{\Lambda}$ is a diagonal matrix of the leading $m$ eigenvalues of $\mathbf{K}$ and $\mathbf{\Phi}$ is the matrix of the corresponding orthonormal eigenvectors (Golub and Van Loan, 1996; Rasmussen and Williams, 2006, Ch. 8). Yet, as computing the eigendecomposition is an $\mathcal{O}(n^3)$ operation, this provides no remedy as such.

In this work we propose a novel method for obtaining approximate eigendecompositions of covariance functions in terms of an eigenfunction expansion of the Laplace operator in a compact subset of $\mathbb{R}^d$. The method is based on interpreting the covariance function as the kernel of a pseudo-differential operator (Shubin, 1987) and approximating it using Hilbert space methods (Courant and Hilbert, 2008; Showalter, 2010). This results in a reduced-rank approximation for the covariance function, where the basis functions are independent of the covariance functions and its parameters. We also show that the approximation converges to the exact solution in well-defined conditions, analyze its convergence rate, and provide theoretical and experimental comparisons to existing state-of-the-art methods. This path has not been explored in GP regression context before, although the approach is related to the Fourier feature methods (Hensman et al, 2018) and stochastic partial differential equation based methods recently introduced to spatial statistics and GP regression (Lindgren et al, 2011; Särkkä and Hartikainen, 2012; Särkkä

et al, 2013) as well as to classical works in the spectral representations of stochastic processes (Loève, 1963; Van Trees, 1968; Adler, 1981; Cramér and Leadbetter, 2013) and spline interpolation (Wahba, 1978; Kimeldorf and Wahba, 1970; Wahba, 1990). Recently, the scalable eigendecomposition approach has also been tackled by various structure exploiting methods (building on the work by Wilson and Nickisch, 2015) and extended to methods exploiting GPU computations.

This paper is structured as follows: In Section 2 we derive the approximative series expansion of the covariance functions. Section 3 is dedicated to applying the approximation scheme to GP regression and providing details of the computational benefits. We provide a detailed analysis of the convergence of the method in Section 4. Section 5 and 6 provide comparisons to existing methods, the former from a more theoretical point of view, whereas the latter contains examples and comparative evaluation on several datasets. Finally the properties of the method are summarized and discussed in Section 7.

## 2 Approximating the Covariance Function

In this section, we start by stating the assumptions and properties of the class of covariance functions that we are considering, and show how a homogenous covariance function can be considered as a pseudo-differential operator constructed as a series of Laplace operators. Then we show how the pseudo-differential operators can be approximated with Hilbert space methods on compact subsets of $\mathbb{R}^d$ or via inner products with integrable weight functions, and discuss connections to Sturm–Liouville theory.

### 2.1 Spectral Densities of Homogeneous and Isotropic Gaussian Processes

In this work it is assumed that the covariance function is homogeneous (stationary), which means that the covariance function $k(\mathbf{x}, \mathbf{x}')$ is actually a function of $\mathbf{r} = \mathbf{x} - \mathbf{x}'$ only. This means that the covariance structure of the model function $f(\mathbf{x})$ is the same regardless of the absolute position in the input space (*cf.* Rasmussen and Williams, 2006, Ch. 4). In this case the covariance function can be equivalently represented in terms of the spectral density. This results from the *Bochner's theorem* (see, *e.g.*, Akhiezer and Glazman, 1993; Da Prato and Zabczyk, 1992) which states that a bounded continuous positive definite function $k(\mathbf{r})$ can be represented

as

$$k(\mathbf{r}) = \frac{1}{(2\pi)^d} \int \exp\left(\mathrm{i}\,\omega^\mathsf{T}\mathbf{r}\right) \mu(\mathrm{d}\omega), \tag{2}$$

where $\mu$ is a positive measure.

If the measure $\mu(\omega)$ has a density, it is called the *spectral density* $S(\omega)$ corresponding to the covariance function $k(\mathbf{r})$. This gives rise to the Fourier duality of covariance and spectral density, which is known as the *Wiener–Khintchin theorem* (Rasmussen and Williams, 2006, Ch. 4), giving the identities

$$k(\mathbf{r}) = \frac{1}{(2\pi)^d} \int S(\omega)\,e^{\mathrm{i}\,\omega^\mathsf{T}\mathbf{r}}\,\mathrm{d}\omega,$$
$$S(\omega) = \int k(\mathbf{r})\,e^{-\mathrm{i}\,\omega^\mathsf{T}\mathbf{r}}\,\mathrm{d}\mathbf{r}. \tag{3}$$

From these identities it is easy to see that if the covariance function is *isotropic*, that is, it only depends on the Euclidean norm $\|\mathbf{r}\|$ such that $k(\mathbf{r}) \triangleq k(\|\mathbf{r}\|)$, then the spectral density will also only depend on the norm of $\omega$ such that we can write $S(\omega) \triangleq S(\|\omega\|)$. In the following we assume that the considered covariance functions are indeed isotropic, but the approach can be generalized to more general homogenous covariance functions.

## 2.2 The Covariance Operator As a Pseudo-Differential Operator

Associated to each covariance function $k(\mathbf{x}, \mathbf{x}')$ we can also define a covariance operator $\mathcal{K}$ as follows:

$$\mathcal{K}\,\phi = \int k(\cdot, \mathbf{x}')\,\phi(\mathbf{x}')\,\mathrm{d}\mathbf{x}'. \tag{4}$$

Note that because the covariance function is homogeneous, this can also be written as a convolution. As we show in the next section, this interpretation allows us to approximate the covariance operator using Hilbert space methods which are typically used for approximating differential and pseudo-differential operators in the context of partial differential equations (Showalter, 2010). When the covariance function is homogenous, the corresponding operator will be translation invariant thus allowing for Fourier-representation as a transfer function. This transfer function is just the spectral density of the Gaussian process.

Consider an isotropic covariance function $k(\mathbf{x}, \mathbf{x}') \triangleq k(\|\mathbf{r}\|)$ (recall that $\|\cdot\|$ denotes the Euclidean norm). The spectral density of the Gaussian process and thus the transfer function corresponding to the covariance operator will now have the form $S(\|\omega\|)$. We can formally write it as a function of $\|\omega\|^2$ such that

$$S(\|\omega\|) = \psi(\|\omega\|^2). \tag{5}$$

Assume that the spectral density $S(\cdot)$ and hence $\psi(\cdot)$ have the following polynomial expansion:

$$\psi(\|\omega\|^2) = a_0 + a_1\|\omega\|^2 + a_2(\|\omega\|^2)^2 + a_3(\|\omega\|^2)^3 + \cdots. \tag{6}$$

This can be ensured, for example, by requiring that $\psi(\cdot)$ is an analytic function. Thus we also have

$$S(\|\omega\|) = a_0 + a_1\|\omega\|^2 + a_2(\|\omega\|^2)^2 + a_3(\|\omega\|^2)^3 + \cdots. \tag{7}$$

Recall that the transfer function corresponding to the Laplace operator $\nabla^2$ is $-\|\omega\|^2$ in the sense that for a regular enough function $f$ we have

$$\mathscr{F}[\nabla^2 f](\omega) = -\|\omega\|^2 \mathscr{F}[f](\omega), \tag{8}$$

where $\mathscr{F}[\cdot]$ denotes the Fourier transform of its argument. If we take the inverse Fourier transform of (7), we get the following representation for the covariance operator $\mathcal{K}$, which defines a pseudo-differential operator (Shubin, 1987) as a formal series of Laplace operators:

$$\mathcal{K} = a_0 + a_1(-\nabla^2) + a_2(-\nabla^2)^2 + a_3(-\nabla^2)^3 + \cdots. \tag{9}$$

In the next section we will use this representation to form a series expansion approximation for the covariance function.

## 2.3 Hilbert-Space Approximation of the Covariance Operator

We will now form a Hilbert-space approximation for the pseudo-differential operator defined by (9). Let $\Omega \subset \mathbb{R}^d$ be a compact set, and consider the eigenvalue problem for the Laplace operators with Dirichlet boundary conditions (we could use other boundary conditions as well):

$$\begin{cases} -\nabla^2 \phi_j(\mathbf{x}) = \lambda_j\,\phi_j(\mathbf{x}), & \mathbf{x} \in \Omega, \\ \phi_j(\mathbf{x}) = 0, & \mathbf{x} \in \partial\Omega. \end{cases} \tag{10}$$

Let us now assume that we have selected $\partial\Omega$ to be sufficiently smooth, for example, a hypercube or hypersphere, so that the eigenfunctions and eigenvalues exist. Because $-\nabla^2$ is a positive definite Hermitian operator, the set of eigenfunctions $\phi_j(\cdot)$ is orthonormal with respect to the inner product

$$\langle f, g \rangle = \int_\Omega f(\mathbf{x})\,g(\mathbf{x})\,\mathrm{d}\mathbf{x} \tag{11}$$

that is,

$$\int_\Omega \phi_i(\mathbf{x})\,\phi_j(\mathbf{x})\,\mathrm{d}\mathbf{x} = \delta_{ij}, \tag{12}$$

and all the eigenvalues $\lambda_j$ are real and positive. The negative Laplace operator can then be assigned the formal kernel

$$l(\mathbf{x}, \mathbf{x}') = \sum_j \lambda_j \, \phi_j(\mathbf{x}) \, \phi_j(\mathbf{x}') \qquad (13)$$

in the sense that

$$-\nabla^2 f(\mathbf{x}) = \int l(\mathbf{x}, \mathbf{x}') \, f(\mathbf{x}') \, \mathrm{d}\mathbf{x}', \qquad (14)$$

for sufficiently (weakly) differentiable functions $f$ in the domain $\Omega$ assuming Dirichlet boundary conditions.

If we consider the formal powers of this representation, due to orthonormality of the basis, we can write the arbitrary operator power $s = 1, 2, \ldots$ of the kernel as

$$l^s(\mathbf{x}, \mathbf{x}') = \sum_j \lambda_j^s \, \phi_j(\mathbf{x}) \, \phi_j(\mathbf{x}'). \qquad (15)$$

This is again to be interpreted to mean that

$$(-\nabla^2)^s f(\mathbf{x}) = \int l^s(\mathbf{x}, \mathbf{x}') \, f(\mathbf{x}') \, \mathrm{d}\mathbf{x}', \qquad (16)$$

for regular enough functions $f$ and in the current domain with the assumed boundary conditions.

This implies that on the domain $\Omega$, assuming the boundary conditions, we also have

$$\left[ a_0 + a_1(-\nabla^2) + a_2(-\nabla^2)^2 + \cdots \right] f(\mathbf{x})$$
$$= \int \left[ a_0 + a_1 \, l^1(\mathbf{x}, \mathbf{x}') + a_2 \, l^2(\mathbf{x}, \mathbf{x}') + \cdots \right] f(\mathbf{x}') \, \mathrm{d}\mathbf{x}'. \qquad (17)$$

The left hand side is just $\mathcal{K} f$ via (9), on the domain with the boundary conditions, and thus by comparing to (4) and using (15) we can conclude that

$$k(\mathbf{x}, \mathbf{x}') \approx a_0 + a_1 \, l^1(\mathbf{x}, \mathbf{x}') + a_2 \, l^2(\mathbf{x}, \mathbf{x}') + \cdots$$
$$= \sum_j \left[ a_0 + a_1 \lambda_j^1 + a_2 \lambda_j^2 + \cdots \right] \phi_j(\mathbf{x}) \, \phi_j(\mathbf{x}'), \qquad (18)$$

which is only an approximation to the covariance function due to restriction of the domain to $\Omega$ and the boundary conditions. By letting $\|\omega\|^2 = \lambda_j$ in (7) we now obtain

$$S(\sqrt{\lambda_j}) = a_0 + a_1 \lambda_j^1 + a_2 \lambda_j^2 + \cdots \qquad (19)$$

and substituting this into (18) then leads to the approximation

$$\boxed{k(\mathbf{x}, \mathbf{x}') \approx \sum_j S(\sqrt{\lambda_j}) \, \phi_j(\mathbf{x}) \, \phi_j(\mathbf{x}'),} \qquad (20)$$

where $S(\cdot)$ is the spectral density of the covariance function, $\lambda_j$ is the $j$th eigenvalue and $\phi_j(\cdot)$ the eigenfunction of the Laplace operator in a given domain. These expressions tend to be simple closed-form expressions.

The right hand side of (20) is very easy to evaluate, because it corresponds to evaluating the spectral density at the square roots of the eigenvalues and multiplying them with the eigenfunctions of the Laplace operator. Because the eigenvalues of the Laplace operator are monotonically increasing with $j$ and for bounded covariance functions the spectral density goes to zero fast with higher frequencies, we can expect to obtain a good approximation of the right hand side by retaining only a finite number of terms in the series. However, even with an infinite number of terms this is only an approximation, because we assumed a compact domain with boundary conditions. The approximation can be, though, expected to be good at the input values which are not near the boundary of $\Omega$, where the Laplacian was taken to be zero.

As an example, Figure 1 shows Matérn covariance functions of various degrees of smoothness $\nu$ (see, e.g., Rasmussen and Williams, 2006, Ch. 4) and approximations for different numbers of basis functions in the approximation. The basis consists of the eigenfunctions of the Laplacian in (10) with $\Omega = [-L, L]$ which gives the eigenfunctions $\phi_j(x) = L^{-1/2} \sin(\pi j(x + L)/(2L))$ and the eigenvalues $\lambda_j = (\pi \, j/(2L))^2$. In the figure we have set $L = 1$ and $\ell = 0.1$. For the squared exponential the approximation is indistinguishable from the exact curve already at $m = 12$, whereas the less smooth functions require more terms.

### 2.4 Inner Product Point of View

Instead of considering a compact bounded set $\Omega$, we can consider the same approximation in terms of an inner product defined by an input density (Williams and Seeger, 2000). Let the inner product be defined as

$$\langle f, g \rangle = \int f(\mathbf{x}) \, g(\mathbf{x}) \, w(\mathbf{x}) \, \mathrm{d}\mathbf{x}, \qquad (21)$$

where $w(\mathbf{x})$ is some positive weight function such that $\int w(\mathbf{x}) \, \mathrm{d}\mathbf{x} < \infty$. In terms of this inner product, we define the operator

$$\mathcal{K} f = \int k(\cdot, \mathbf{x}) \, f(\mathbf{x}) \, w(\mathbf{x}) \, \mathrm{d}\mathbf{x}. \qquad (22)$$

This operator is self-adjoint with respect to the inner product, $\langle \mathcal{K} f, g \rangle = \langle f, \mathcal{K} g \rangle$, and according to the spectral theorem there exists an orthonormal set of basis functions and positive constants, $\{\varphi_j(\mathbf{x}), \gamma_j \mid j =$
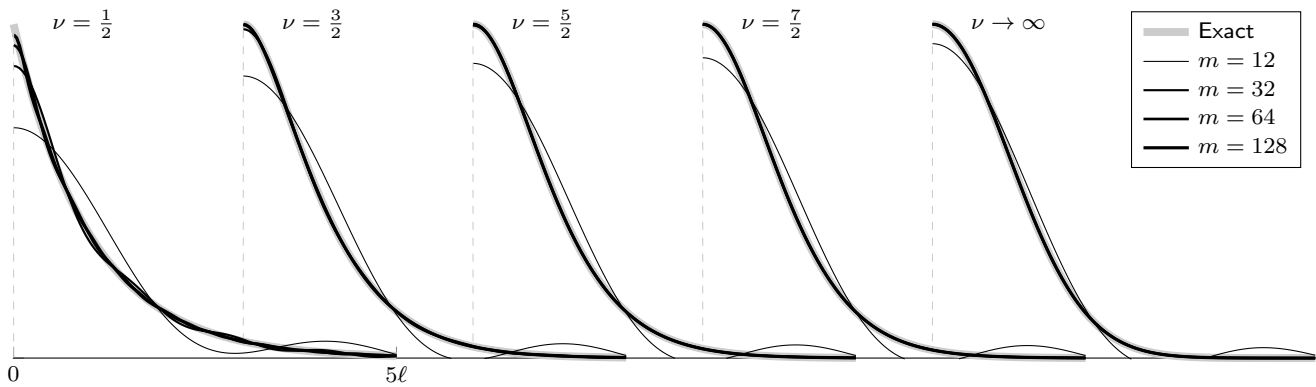
Fig. 1: Approximations to covariance functions of the Matérn class of various degrees of smoothness; $\nu = 1/2$ corresponds to the exponential Ornstein–Uhlenbeck covariance function, and $\nu \to \infty$ to the squared exponential (exponentiated quadratic) covariance function. Approximations are shown for 12, 32, 64, and 128 eigenfunctions.

$1, 2, \ldots\}$, that satisfies the eigenvalue equation

$$(\mathcal{K}\varphi_j)(\mathbf{x}) = \gamma_j\,\varphi_j(\mathbf{x}). \tag{23}$$

Thus $k(\mathbf{x}, \mathbf{x}')$ has the series expansion

$$k(\mathbf{x}, \mathbf{x}') = \sum_j \gamma_j\,\varphi_j(\mathbf{x})\,\varphi_j(\mathbf{x}'). \tag{24}$$

Similarly, we also have the *Karhunen–Loeve expansion* for a sample function $f(\mathbf{x})$ with zero mean and the above covariance function:

$$f(\mathbf{x}) = \sum_j f_j\,\varphi_j(\mathbf{x}), \tag{25}$$

where $f_j$s are independent zero mean Gaussian random variables with variances $\gamma_j$ (see, *e.g.*, Lenk, 1991).

For the negative Laplacian the corresponding definition is

$$\mathcal{D}f = -\nabla^2[f\,w], \tag{26}$$

which implies

$$\langle \mathcal{D}f, g \rangle = -\int f(\mathbf{x})\,w(\mathbf{x})\nabla^2[g(\mathbf{x})\,w(\mathbf{x})]\,\mathrm{d}\mathbf{x}, \tag{27}$$

and the operator defined by (26) can be seen to be self-adjoint. Again, there exists an orthonormal basis $\{\phi_j(\mathbf{x}) \mid j = 1, 2, \ldots\}$ and positive eigenvalues $\lambda_j$ which satisfy the eigenvalue equation

$$(\mathcal{D}\,\phi_j)(\mathbf{x}) = \lambda_j\,\phi_j(\mathbf{x}). \tag{28}$$

Thus the kernel of $\mathcal{D}$ has a series expansion similar to Equation (13) and thus an approximation can be given in the same form as in Equation (20). In this case the approximation error comes from approximating the Laplace operator with the more smooth operator,

$$\nabla^2 f \approx \nabla^2[f\,w], \tag{29}$$

which is closely related to assumption of an input density $w(\mathbf{x})$ for the Gaussian process. However, when the weight function $w(\cdot)$ is close to constant in the area where the inputs points are located, the approximation is accurate.

## 2.5 Connection to Sturm–Liouville Theory

The presented methodology is also related to the Sturm–Liouville theory arising in the theory of partial differential equations (Courant and Hilbert, 2008). When the input $x$ is scalar, the eigenvalue problem in Equation (23) can be written in Sturm–Liouville form as follows:

$$-\frac{\mathrm{d}}{\mathrm{d}x}\left[w^2(x)\,\frac{\mathrm{d}\phi_j(x)}{\mathrm{d}x}\right] - w(x)\,\frac{\mathrm{d}^2 w(x)}{\mathrm{d}x^2}\,\phi_j(x)$$
$$= \lambda_j\,w(x)\,\phi_j(x). \tag{30}$$

The above equation can be solved for $\phi_j(x)$ and $\lambda_j$ using numerical methods for Sturm–Liouville equations. Also note that if we select $w(x) = 1$ in a finite set, we obtain the equation $-\mathrm{d}^2/\mathrm{d}x^2\,\phi_j(x) = \lambda_j\,\phi_j(x)$ which is compatible with the results in the previous section.

We consider the case where $\mathbf{x} \in \mathbb{R}^d$ and $w(\mathbf{x})$ is symmetric around the origin and thus is only a function of the norm $r = \|\mathbf{x}\|$ (*i.e.* has the form $w(r)$). The Laplacian in spherical coordinates is

$$\nabla^2 f = \frac{1}{r^{d-1}}\frac{\partial}{\partial r}\left(r^{d-1}\frac{\partial f}{\partial r}\right) + \frac{1}{r^2}\,\Delta_{S^{d-1}}f, \tag{31}$$

where $\Delta_{S^{d-1}}$ is the Laplace–Beltrami operator on $S^{d-1}$. Let us assume that $\phi_j(r, \xi) = h_j(r)\,g(\xi)$, where $\xi$ denotes the angular variables. After some algebra, writing

(a) $\nu = \frac{1}{2}$ and $\ell = 0.5$       (b) $\nu = \frac{3}{2}$ and $\ell = 0.5$       (c) $\nu \to \infty$ and $\ell = 0.5$
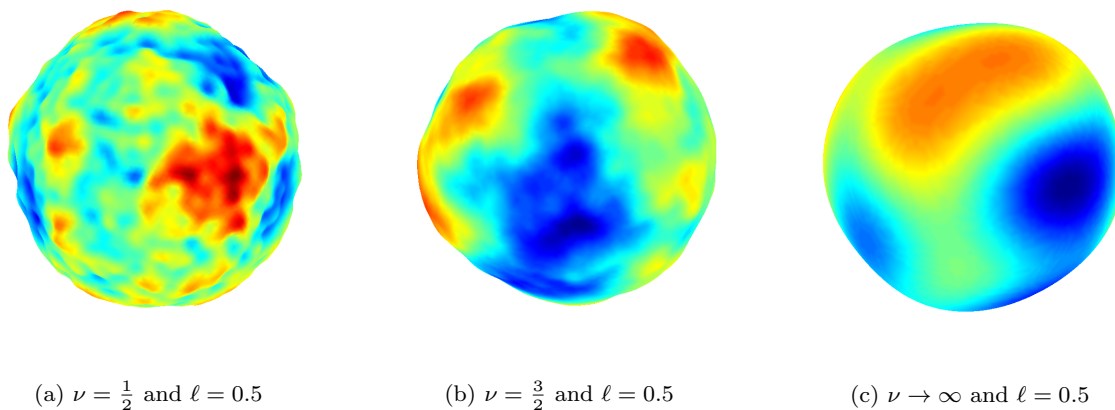
Fig. 2: Approximate random draws of Gaussian processes with the Matérn covariance function on the hull of a unit sphere. The color scale and radius follow the process.

the equations into Sturm–Liouville form yields for the radial part

$$
-\frac{\mathrm{d}}{\mathrm{d}r}\left(w^2(r)\,r\,\frac{\mathrm{d}h_j(r)}{\mathrm{d}r}\right) \\
-\left(\frac{\mathrm{d}w(r)}{\mathrm{d}r}\,w(r) + \frac{\mathrm{d}^2w(r)}{\mathrm{d}r^2}\,w(r)\,r\right)h_j(r) \\
= \lambda_j\,w(r)\,r\,h_j(r), \quad (32)
$$

and $\Delta_{S^{d-1}}g(\xi) = 0$ for the angular part. The solutions to the angular part are the Laplace's spherical harmonics. Note that if we assume that we have $w(r) = 1$ on some area of finite radius, the first equation becomes (when $d > 1$):

$$
r^2\,\frac{\mathrm{d}^2h_j(r)}{\mathrm{d}r^2} + r\,\frac{\mathrm{d}h_j(r)}{\mathrm{d}r} + r^2\,\lambda_j\,h_j(r) = 0. \quad (33)
$$

Figure 2 shows example Gaussian random field draws on a unit sphere, where the basis functions are the Laplace spherical harmonics and the covariance functions of the Matérn class with different degrees of smoothness $\nu$. Our approximation is straight-forward to apply in any domain, where the eigendecomposition of the Laplacian can be formed.

## 3 Application of the Method to GP Regression

In this section we show how the approximation (20) can be used in Gaussian process regression. We also write down the expressions needed for hyperparameter learning and discuss the computational requirements of the methods.

### 3.1 Gaussian Process Regression

GP regression is usually formulated as predicting an unknown scalar output $f(\mathbf{x}_*)$ associated with a known input $\mathbf{x}_* \in \mathbb{R}^d$, given a training data set $\mathcal{D} = \{(\mathbf{x}_i, y_i) \mid i = 1, 2, \ldots, n\}$. The model functions $f$ are assumed to be realizations of a Gaussian random process prior and the observations corrupted by Gaussian noise:

$$
\begin{aligned}
f &\sim \mathcal{GP}(0, k(\mathbf{x}, \mathbf{x}')) \\
y_i &= f(\mathbf{x}_i) + \varepsilon_i,
\end{aligned} \quad (34)
$$

where $\varepsilon_i \sim \mathcal{N}(0, \sigma_{\mathrm{n}}^2)$. For notational simplicity the functions in the above model are *a priori* zero mean and the measurement errors are independent Gaussian, but the results of this paper can be easily generalized to arbitrary mean functions and dependent Gaussian errors. The direct solution to the GP regression problem (34) gives the predictions $p(f(\mathbf{x}_*) \mid \mathcal{D}) = \mathcal{N}(f(\mathbf{x}_*) \mid \mathbb{E}[f(\mathbf{x}_*)], \mathbb{V}[f(\mathbf{x}_*)])$. The conditional mean and variance can be computed in closed-form as (see, *e.g.*, Rasmussen and Williams, 2006, p. 17)

$$
\begin{aligned}
\mathbb{E}[f(\mathbf{x}_*)] &= \mathbf{k}_*^\mathsf{T}(\mathbf{K} + \sigma_{\mathrm{n}}^2\mathbf{I})^{-1}\mathbf{y}, \\
\mathbb{V}[f(\mathbf{x}_*)] &= k(\mathbf{x}_*, \mathbf{x}_*) - \mathbf{k}_*^\mathsf{T}(\mathbf{K} + \sigma_{\mathrm{n}}^2\mathbf{I})^{-1}\mathbf{k}_*,
\end{aligned} \quad (35)
$$

where $\mathbf{K}_{ij} = k(\mathbf{x}_i, \mathbf{x}_j)$, $\mathbf{k}_*$ is an $n$-dimensional vector with the $i$th entry being $k(\mathbf{x}_*, \mathbf{x}_i)$, and $\mathbf{y}$ is a vector of the $n$ observations.

In order to avoid the $n \times n$ matrix inversion in (35), we use the approximation scheme presented in the previous section and project the process to a truncated set of $m$ basis functions of the Laplacian as given in Equation (20) such that

$$
f(\mathbf{x}) \approx \sum_{j=1}^{m} f_j\,\phi_j(\mathbf{x}), \quad (36)
$$

where $f_j \sim \mathcal{N}(0, S(\sqrt{\lambda_j}))$. We can then form an approximate eigendecomposition of the matrix $\mathbf{K} \approx \boldsymbol{\Phi}\boldsymbol{\Lambda}\boldsymbol{\Phi}^{\mathsf{T}}$, where $\boldsymbol{\Lambda}$ is a diagonal matrix of the leading $m$ approximate eigenvalues such that $\boldsymbol{\Lambda}_{jj} = S(\sqrt{\lambda_j}), j = 1, 2, \ldots, m$. Here $S(\cdot)$ is the spectral density of the Gaussian process and $\lambda_j$ the $j$th eigenvalue of the Laplace operator. The corresponding eigenvectors in the decomposition are given by the eigenvectors $\phi_j(\mathbf{x})$ of the Laplacian such that $\boldsymbol{\Phi}_{ij} = \phi_j(\mathbf{x}_i)$.

Using the matrix inversion lemma we rewrite (35) as follows:

$$\begin{aligned}
\mathbb{E}[f_*] &\approx \phi_*^{\mathsf{T}}(\boldsymbol{\Phi}^{\mathsf{T}}\boldsymbol{\Phi} + \sigma_{\mathrm{n}}^2\boldsymbol{\Lambda}^{-1})^{-1}\boldsymbol{\Phi}^{\mathsf{T}}\mathbf{y}, \\
\mathbb{V}[f_*] &\approx \sigma_{\mathrm{n}}^2\phi_*^{\mathsf{T}}(\boldsymbol{\Phi}^{\mathsf{T}}\boldsymbol{\Phi} + \sigma_{\mathrm{n}}^2\boldsymbol{\Lambda}^{-1})^{-1}\phi_*,
\end{aligned} \quad (37)$$

where $\phi_*$ is an $m$-dimensional vector with the $j$th entry being $\phi_j(\mathbf{x}_*)$. Thus, when the size of the training set is higher than the number of required basis functions $n > m$, the use of this approximation is advantageous.

### 3.2 Learning the Hyperparameters

A common way to learn the hyperparameters $\theta$ of the covariance function (suppressed earlier in the notation for brevity) and the noise variance $\sigma_{\mathrm{n}}^2$ is by maximizing the marginal likelihood function (Rasmussen and Williams, 2006; Quiñonero-Candela and Rasmussen, 2005a). Let $\mathbf{Q} = \mathbf{K} + \sigma_{\mathrm{n}}^2\mathbf{I}$ for the full model, then the negative log marginal likelihood and its derivatives are

$$\mathcal{L} = \frac{1}{2}\log|\mathbf{Q}| + \frac{1}{2}\mathbf{y}^{\mathsf{T}}\mathbf{Q}^{-1}\mathbf{y} + \frac{n}{2}\log(2\pi), \quad (38)$$

$$\frac{\partial\mathcal{L}}{\partial\theta_k} = \frac{1}{2}\mathrm{Tr}\left(\mathbf{Q}^{-1}\frac{\partial\mathbf{Q}}{\partial\theta_k}\right) - \frac{1}{2}\mathbf{y}^{\mathsf{T}}\mathbf{Q}^{-1}\frac{\partial\mathbf{Q}}{\partial\theta_k}\mathbf{Q}^{-1}\mathbf{y}, \quad (39)$$

$$\frac{\partial\mathcal{L}}{\partial\sigma_{\mathrm{n}}^2} = \frac{1}{2}\mathrm{Tr}\left(\mathbf{Q}^{-1}\right) - \frac{1}{2}\mathbf{y}^{\mathsf{T}}\mathbf{Q}^{-1}\mathbf{Q}^{-1}\mathbf{y}, \quad (40)$$

and they can be combined with a conjugate gradient optimizer. The problem in this case is the inversion of $\mathbf{Q}$, which is an $n \times n$ matrix. And thus each step of running the optimizer is $\mathcal{O}(n^3)$. For our approximation scheme, let $\tilde{\mathbf{Q}} = \boldsymbol{\Phi}\boldsymbol{\Lambda}\boldsymbol{\Phi}^{\mathsf{T}} + \sigma_{\mathrm{n}}^2\mathbf{I}$. Now replacing $\mathbf{Q}$ with $\tilde{\mathbf{Q}}$ in the above expressions gives us the following:

$$\tilde{\mathcal{L}} = \frac{1}{2}\log|\tilde{\mathbf{Q}}| + \frac{1}{2}\mathbf{y}^{\mathsf{T}}\tilde{\mathbf{Q}}^{-1}\mathbf{y} + \frac{n}{2}\log(2\pi), \quad (41)$$

$$\frac{\partial\tilde{\mathcal{L}}}{\partial\theta_k} = \frac{1}{2}\frac{\partial\log|\tilde{\mathbf{Q}}|}{\partial\theta_k} + \frac{1}{2}\frac{\partial\mathbf{y}^{\mathsf{T}}\tilde{\mathbf{Q}}^{-1}\mathbf{y}}{\partial\theta_k}, \quad (42)$$

$$\frac{\partial\tilde{\mathcal{L}}}{\partial\sigma_{\mathrm{n}}^2} = \frac{1}{2}\frac{\partial\log|\tilde{\mathbf{Q}}|}{\partial\sigma_{\mathrm{n}}^2} + \frac{1}{2}\frac{\partial\mathbf{y}^{\mathsf{T}}\tilde{\mathbf{Q}}^{-1}\mathbf{y}}{\partial\sigma_{\mathrm{n}}^2}, \quad (43)$$

where for the terms involving $\log|\tilde{\mathbf{Q}}|$:

$$\begin{aligned}
\log|\tilde{\mathbf{Q}}| = {}&(n-m)\log\sigma_{\mathrm{n}}^2 + \log|\mathbf{Z}| \\
&+ \sum_{j=1}^m \log S(\sqrt{\lambda_j}),
\end{aligned} \quad (44)$$

$$\begin{aligned}
\frac{\partial\log|\tilde{\mathbf{Q}}|}{\partial\theta_k} = {}&\sum_{j=1}^m S(\sqrt{\lambda_j})^{-1}\frac{\partial S(\sqrt{\lambda_j})}{\partial\theta_k} \\
&- \sigma_{\mathrm{n}}^2\mathrm{Tr}\left(\mathbf{Z}^{-1}\boldsymbol{\Lambda}^{-2}\frac{\partial\boldsymbol{\Lambda}}{\partial\theta_k}\right),
\end{aligned} \quad (45)$$

$$\frac{\partial\log|\tilde{\mathbf{Q}}|}{\partial\sigma_{\mathrm{n}}^2} = \frac{n-m}{\sigma_{\mathrm{n}}^2} + \mathrm{Tr}\left(\mathbf{Z}^{-1}\boldsymbol{\Lambda}^{-1}\right), \quad (46)$$

and for the terms involving $\tilde{\mathbf{Q}}^{-1}$:

$$\mathbf{y}^{\mathsf{T}}\tilde{\mathbf{Q}}^{-1}\mathbf{y} = \frac{1}{\sigma_{\mathrm{n}}^2}\left(\mathbf{y}^{\mathsf{T}}\mathbf{y} - \mathbf{y}^{\mathsf{T}}\boldsymbol{\Phi}\mathbf{Z}^{-1}\boldsymbol{\Phi}^{\mathsf{T}}\mathbf{y}\right), \quad (47)$$

$$\frac{\partial\mathbf{y}^{\mathsf{T}}\tilde{\mathbf{Q}}^{-1}\mathbf{y}}{\partial\theta_k} = -\mathbf{y}^{\mathsf{T}}\boldsymbol{\Phi}\mathbf{Z}^{-1}\left(\boldsymbol{\Lambda}^{-2}\frac{\partial\boldsymbol{\Lambda}}{\partial\theta_k}\right)\mathbf{Z}^{-1}\boldsymbol{\Phi}^{\mathsf{T}}\mathbf{y}, \quad (48)$$

$$\frac{\partial\mathbf{y}^{\mathsf{T}}\tilde{\mathbf{Q}}^{-1}\mathbf{y}}{\partial\sigma_{\mathrm{n}}^2} = \frac{1}{\sigma_{\mathrm{n}}^2}\mathbf{y}^{\mathsf{T}}\boldsymbol{\Phi}\mathbf{Z}^{-1}\boldsymbol{\Lambda}^{-1}\mathbf{Z}^{-1}\boldsymbol{\Phi}^{\mathsf{T}}\mathbf{y} - \frac{1}{\sigma_{\mathrm{n}}^4}\mathbf{y}^{\mathsf{T}}\mathbf{y}, \quad (49)$$

where $\mathbf{Z} = \sigma_{\mathrm{n}}^2\boldsymbol{\Lambda}^{-1} + \boldsymbol{\Phi}^{\mathsf{T}}\boldsymbol{\Phi}$. For efficient implementation, matrix-to-matrix multiplications can be avoided in many cases, and the inversion of $\mathbf{Z}$ can be carried out through Cholesky factorization for numerical stability. This factorization ($\mathbf{L}\mathbf{L}^{\mathsf{T}} = \mathbf{Z}$) can also be used for the term $\log|\mathbf{Z}| = 2\sum_j \log\mathbf{L}_{jj}$, and $\mathrm{Tr}\left(\mathbf{Z}^{-1}\boldsymbol{\Lambda}^{-1}\right) = \sum_j 1/(\mathbf{Z}_{jj}\boldsymbol{\Lambda}_{jj})$ can be evaluated by element-wise multiplication.

Once the marginal likelihood and its derivatives are available, it is also possible to use other methods for parameter inference such as Markov chain Monte Carlo methods (Liu, 2001; Brooks et al, 2011) including Hamiltonian Monte Carlo (HMC, Duane et al, 1987; Neal, 2011) as well as numerous others.

### 3.3 Discussion on the Computational Complexity

As can be noted from Equation (20), the basis functions in the reduced-rank approximation do not depend on the hyperparameters of the covariance function. Thus it is enough to calculate the product $\boldsymbol{\Phi}^{\mathsf{T}}\boldsymbol{\Phi}$ only once, which means that the method has a overall asymptotic computational complexity of $\mathcal{O}(nm^2)$. After this initial cost, evaluating the marginal likelihood and the marginal likelihood gradient is an $\mathcal{O}(m^3)$ operation—which in practice comes from the Cholesky factorization of $\mathbf{Z}$ on each step.

If the number of observations $n$ is so large that storing the $n \times m$ matrix $\boldsymbol{\Phi}$ is not feasible, the computations

of $\boldsymbol{\Phi}^\mathsf{T}\boldsymbol{\Phi}$ can be carried out in blocks. Storing the evaluated eigenfunctions in $\boldsymbol{\Phi}$ is not necessary, because the $\phi_j(\mathbf{x})$ are closed-form expressions that can be evaluated when necessary. In practice, it might be preferable to cache the result of $\boldsymbol{\Phi}^\mathsf{T}\boldsymbol{\Phi}$ (causing a memory requirement scaling as $\mathcal{O}(m^2)$), but this is not required.

The computational complexity of conventional sparse GP approximations typically scale as $\mathcal{O}(nm^2)$ in time for each step of evaluating the marginal likelihood. The scaling in demand of storage is $\mathcal{O}(nm)$. This comes from the inevitable cost of re-evaluating all results involving the basis functions on each step and storing the matrices required for doing this. This applies to all the methods that will be discussed in Section 5, with the exception of SSGP, where the storage demand can be relaxed by re-evaluating the basis functions on demand.

We can also consider the rather restricting, but in certain applications often encountered case, where the measurements are constrained to a regular grid. This causes the product of the orthonormal eigenfunction matrices $\boldsymbol{\Phi}^\mathsf{T}\boldsymbol{\Phi}$ to be diagonal, avoiding the calculation of the matrix inverse altogether. This relates to the FFT-based methods for GP regression (Paciorek, 2007; Fritz et al, 2009), and the projections to the basis functions can be evaluated by fast Fourier transform in $\mathcal{O}(n \log n)$ time complexity.

### 3.4 Inverse Problems and Latent Force Models

We can also use the methodology to models of the form

$$
\begin{aligned}
f(\mathbf{x}) &\sim \mathcal{GP}(0, k(\mathbf{x}, \mathbf{x}')), \\
y_i &= (\mathcal{H}f)(\mathbf{x}_i) + \varepsilon_i,
\end{aligned}
\tag{50}
$$

where $\mathcal{H}$ is a linear operator acting on functions depending on the $\mathbf{x}$ variable. This kind of models appear both in inverse problems literature and machine learning (see, *e.g.*, Tarantola, 2004; Kaipio and Somersalo, 2005; Särkkä, 2011). The Gaussian process regression solution now becomes

$$
\begin{aligned}
\mathbb{E}[f(\mathbf{x}_*)] &= \mathbf{k}_{*h}^\mathsf{T}(\mathbf{K}_h + \sigma_\mathrm{n}^2\mathbf{I})^{-1}\mathbf{y}, \\
\mathbb{V}[f(\mathbf{x}_*)] &= k(\mathbf{x}_*, \mathbf{x}_*) - \mathbf{k}_{*h}^\mathsf{T}(\mathbf{K}_h + \sigma_\mathrm{n}^2\mathbf{I})^{-1}\mathbf{k}_{*h},
\end{aligned}
\tag{51}
$$

where $[\mathbf{K}_h]ij = (\mathcal{H}\,\mathcal{H}'\,k)(\mathbf{x}_i, \mathbf{x}_j)$, the $i$th entry of vector $\mathbf{k}_{*h}$ is $(\mathcal{H}'\,k(\mathbf{x}_*, \cdot))(\mathbf{x}_i)$, and $\mathbf{y}$ is the vector of observations. Here $\mathcal{H}'$ denotes that the operator is applied to the second variable $\mathbf{x}'$ of the argument. With the series expansion (20) we can easily approximate

$$
(\mathcal{H}\,\mathcal{H}'\,k)(\mathbf{x}, \mathbf{x}') \approx \sum_j S(\sqrt{\lambda_j})\,(\mathcal{H}\,\phi_j)(\mathbf{x})\,(\mathcal{H}\,\phi_j)(\mathbf{x}'),
$$

$$
(\mathcal{H}'\,k(\mathbf{x}_*, \cdot))(\mathbf{x}') \approx \sum_j S(\sqrt{\lambda_j})\,\phi_j(\mathbf{x}_*)\,(\mathcal{H}\,\phi_j)(\mathbf{x}').
$$

$$
\tag{52}
$$

After applying the matrix inversion lemma (51) becomes

$$
\begin{aligned}
\mathbb{E}[f_*] &\approx \phi_*^\mathsf{T}(\tilde{\boldsymbol{\Phi}}^\mathsf{T}\tilde{\boldsymbol{\Phi}} + \sigma_\mathrm{n}^2\boldsymbol{\Lambda}^{-1})^{-1}\tilde{\boldsymbol{\Phi}}^\mathsf{T}\mathbf{y}, \\
\mathbb{V}[f_*] &\approx \sigma_\mathrm{n}^2\phi_*^\mathsf{T}(\tilde{\boldsymbol{\Phi}}^\mathsf{T}\tilde{\boldsymbol{\Phi}} + \sigma_\mathrm{n}^2\boldsymbol{\Lambda}^{-1})^{-1}\phi_*,
\end{aligned}
\tag{53}
$$

where $\tilde{\boldsymbol{\Phi}}_{ij} = (\mathcal{H}\phi_j)(\mathbf{x}_i)$ and $\phi_*$ is as defined in (37). The hyperparameter estimation methods discussed in Section 3.2 can also be easily extended to this case.

Another (related) type of model is the following model arising in the context of latent force models (LFM, Álvarez et al, 2013)

$$
\begin{aligned}
f(\mathbf{x}) &\sim \mathcal{GP}(0, k(\mathbf{x}, \mathbf{x}')), \\
\mathcal{L}g &= f, \\
y_i &= g(\mathbf{x}_i) + \varepsilon_i,
\end{aligned}
\tag{54}
$$

where $\mathcal{L}$ is a linear operator. We can now write $\mathcal{H} = \mathcal{L}^{-1}$, where $\mathcal{L}^{-1}$ is the Green's operator associated with the operator $\mathcal{L}$ and hence the model becomes a special case of (50). The approximation to the operator $\mathcal{L}^{-1}$ on the given basis can be easily formed by using, for example, by projecting it onto the basis or by using point collocation. A particularly simple cases arises when the operator itself contains of Laplace operators, for example, when it has the form $\mathcal{L} = \nabla^2$. In that case the projection of the operator becomes diagonal.

## 4 Convergence Analysis

In this section we analyze the convergence of the proposed approximation when the size of the domain $\Omega$ and the number of terms in the series grows to infinity. We start by analyzing a univariate problem in the domain $\Omega = [-L, L]$ and with Dirichlet boundary conditions and then generalize the result to $d$-dimensional cubes $\Omega = [-L_1, L_1] \times \cdots \times [-L_d, L_d]$. Then we analyze the truncation error as function of the number of terms in the series. We also discuss how the analysis could be extended to other types of basis functions.

### 4.1 Univariate Dirichlet Case

In the univariate case, the $m$-term approximation has the form

$$
\widetilde{k}_m(x, x') = \sum_{j=1}^m S(\sqrt{\lambda_j})\,\phi_j(x)\,\phi_j(x'),
\tag{55}
$$

where the eigenfunctions and eigenvalues for $j = 1, 2, \ldots$ are:

$$\phi_j(x) = \frac{1}{\sqrt{L}} \sin\left(\frac{\pi j (x + L)}{2L}\right) \quad \text{and} \quad \lambda_j = \left(\frac{\pi j}{2L}\right)^2. \tag{56}$$

The true covariance function $k(x, x')$ is assumed to be stationary and have a spectral density with the following properties. It is uniformly bounded $S(\omega) = B < \infty$ and has at least one bounded derivative $|S'(\omega)| = D < \infty$ on $\omega > 0$. The following integrals are also assumed to be bounded: $\int_0^\infty S(\omega) \, d\omega = A < \infty$ and $\int_0^\infty |S'(\omega)| \, d\omega = C < \infty$. We also assume that our training and test sets are constrained in the area $[-\widetilde{L}, \widetilde{L}]$, where $\widetilde{L} < L$, and thus we are only interested in the case $x, x' \in [-\widetilde{L}, \widetilde{L}]$. For the purposes of analysis we also assume that $L$ is bounded below by a constant.

The univariate convergence result can be summarized as the following theorem which is proved in Appendix A.2.

**Theorem 1.** *There exists a constant $E$ (independent of $m$, $x$, and $x'$) such that*

$$\left| k(x, x') - \widetilde{k}_m(x, x') \right| \leq \frac{E}{L} + \frac{2}{\pi} \int_{\frac{\pi m}{2L}}^\infty S(\omega) \, d\omega, \tag{57}$$

*which in turn implies that uniformly*

$$\lim_{L \to \infty} \left[ \lim_{m \to \infty} \widetilde{k}_m(x, x') \right] = k(x, x'). \tag{58}$$

**Remark 2.** *Note that we cannot simply exchange the order of the limits in the above theorem. However, the theorem does ensure the convergence of the approximation in the joint limit $m, L \to \infty$ provided that we add terms to the series fast enough such that $m/L \to \infty$. That is, in this limit, the approximation $\widetilde{k}_m(x, x')$ converges uniformly to $k(x, x')$.*

As such, the results above only ensure the convergence of the prior covariance functions. However, it turns out that this also ensures the convergence of the posterior as is summarized in the following corollary.

**Corollary 3.** *Because the Gaussian process regression equations only involve pointwise evaluations of the kernels, it also follows that the posterior mean and covariance functions converge uniformly to the exact solutions in the limit $m, L \to \infty$.*

*Proof.* Analogous to proof of Theorem 2.2 in Särkkä and Piché (2014). □

## 4.2 Multivariate Cartesian Dirichlet Case

In order to generalize the results from the previous section, we turn our attention to a $d$-dimensional inputs space with rectangular domain $\Omega = [-L_1, L_1] \times \cdots \times [-L_d, L_d]$ with Dirichlet boundary conditions. In this case we consider a truncated $m = \hat{m}^d$ term approximation of the form

$$\widetilde{k}_m(\mathbf{x}, \mathbf{x}')$$
$$= \sum_{j_1, \ldots, j_d = 1}^{\hat{m}} S(\sqrt{\lambda_{j_1, \ldots, j_d}}) \, \phi_{j_1, \ldots, j_d}(\mathbf{x}) \, \phi_{j_1, \ldots, j_d}(\mathbf{x}') \tag{59}$$

with the eigenfunctions and eigenvalues

$$\phi_{j_1, \ldots, j_d}(x) = \prod_{k=1}^d \frac{1}{\sqrt{L_k}} \sin\left(\frac{\pi j_k (x_k + L_k)}{2L_k}\right) \tag{60}$$

and

$$\lambda_{j_1, \ldots, j_d} = \sum_{k=1}^d \left(\frac{\pi j_k}{2L_k}\right)^2. \tag{61}$$

The true covariance function $k(\mathbf{x}, \mathbf{x}')$ is assumed to be homogeneous (stationary) and have a spectral density $S(\omega)$ which satisfies the one-dimensional assumptions listed in the previous section in each variable. Furthermore, we assume that the training and test sets are contained in the $d$-dimensional cube $[-\widetilde{L}, \widetilde{L}]^d$ and that $L_k$s are bounded from below.

The following result for this $d$-dimensional case is proved in Appendix A.3.

**Theorem 4.** *There exists a constant $E$ (independent of $m$, $d$, $\mathbf{x}$, and $\mathbf{x}'$) such that*

$$\left| k(\mathbf{x}, \mathbf{x}') - \widetilde{k}_m(\mathbf{x}, \mathbf{x}') \right| \leq \frac{E \, d}{L} + \frac{1}{\pi^d} \int_{\|\omega\| \geq \frac{\pi \hat{m}}{2L}} S(\omega) \, d\omega, \tag{62}$$

*where $L = \min_k L_k$, which in turn implies that uniformly*

$$\lim_{L_1, \ldots, L_d \to \infty} \left[ \lim_{m \to \infty} \widetilde{k}_m(\mathbf{x}, \mathbf{x}') \right] = k(\mathbf{x}, \mathbf{x}'). \tag{63}$$

**Remark 5.** *Analogously as in the one-dimensional case we cannot simply exchange the order of the limits above. Furthermore, we need to add terms fast enough so that $\hat{m}/L_k \to \infty$ when $m, L_1, \ldots, L_d \to \infty$.*

**Corollary 6.** *As in the one-dimensional case, the uniform convergence of the prior covariance function also implies uniform convergence of the posterior mean and covariance in the limit $m, L_1, \ldots, L_d \to \infty$.*

4.3 Scaling of Error with Increasing $\hat{m}$

Using the Dirichlet eigenfunction basis, we can also investigate the truncation error with an increasing number of series expansion terms $m = \hat{m}^d$. If we take a look at the bound in Theorem 4, we can see that it has the form

$$\frac{E\,d}{L} + \frac{1}{\pi^d} \int_{\|\omega\| \geq \frac{\pi\,\hat{m}}{2L}} S(\omega)\,\mathrm{d}\omega, \tag{64}$$

where the first term is independent of $\hat{m}$ and is a linear function of $d$. The latter term in turn depends on $\hat{m}$ and in that sense defines the scaling of error in the number of series terms.

It is worth noting that due to Remarks 17 and 20 we could actually tighten the bound by introducing $\hat{m}$-dependence to $E$, but it does not affect the order of scaling, because the dependence on the dimensionality in that term is linear. Furthermore, the latter term actually depends on the ratio $\hat{m}/L$ and hence there is a coupling between the number of terms and the size of the domain $L$. However, we can still get idea of the convergence speed by fixing $L$.

Let us start by considering the case when $S(\|\omega\|)$ is bounded by a reciprocal of a polynomial which is the case, for example, for the Matérn covariance function. We get the following theorem.

**Theorem 7.** *Assume that where exists a constant $D$ such that $S(\|\omega\|) \leq \frac{D}{\|\omega\|^{d+a}}$ for some $a > 0$. Then we have*

$$\int_{\|\omega\| \geq \frac{\pi\,\hat{m}}{2L}} S(\|\omega\|)\,\mathrm{d}\omega \leq \frac{D'}{m^{a/d}}, \tag{65}$$

*where $m = \hat{m}^d$ for some constant $D'$ (which depends on $L$ and $d$).*

*Proof.* First recall that

$$\int_{\|\omega\| \geq \frac{\pi\,\hat{m}}{2L}} S(\|\omega\|)\,\mathrm{d}\omega = \frac{2\pi^{d/2}}{\Gamma(d/2)} \int_{\frac{\pi\,\hat{m}}{2L}}^{\infty} S(r)\,r^{d-1}\,\mathrm{d}r, \tag{66}$$

where $\Gamma(\cdot)$ is the gamma function, and hence to analyze the scaling as function of $m$, it is enough to investigate the scaling of the term $\int_{\frac{\pi\,\hat{m}}{2L}}^{\infty} S(r)\,r^{d-1}\,dr$. We now get

$$\begin{aligned}
\int_{\frac{\pi\,\hat{m}}{2L}}^{\infty} S(r)\,r^{d-1}\,dr &\leq \int_{\frac{\pi\,\hat{m}}{2L}}^{\infty} \frac{D}{r^{a+1}}\,dr \\
&= \left(\frac{1}{a}\right)\left(\frac{2L}{\pi\,\hat{m}}\right)^{a} \\
&= \underbrace{\left(\frac{(2L)^a}{\pi^a\,a}\right)}_{D'}\left(\frac{1}{m^{a/d}}\right),
\end{aligned} \tag{67}$$

where we have recalled that $m = \hat{m}^d$.                     $\square$

The result in the above theorem tells that by selecting an appropriate differentiation order for the covariance function, we can make the convergence speed arbitrarily large. In particular, if we select $a = d/2$, we get the Monte Carlo rate, and with $a = d$, we get a convergence rate of $\sim 1/m$.

In order to analyze the squared exponential covariance function with spectral density

$$S(\omega) = \prod_{i=1}^{d} \left[s^2\,\sqrt{2\pi}\,\ell\,\exp\left(-\frac{\ell^2\,\omega_i^2}{2}\right)\right], \tag{68}$$

we recall that the integral $\int_{\|\omega\| \geq \frac{\pi\,\hat{m}}{2L}} S(\|\omega\|)\,\mathrm{d}\omega$ was actually used for bounding a more tight bound $\int_{\frac{\pi\,\hat{m}}{2L_1}}^{\infty} \cdots \int_{\frac{\pi\,\hat{m}}{2L_d}}^{\infty} S(\omega_1,\ldots,\omega_d)\,\mathrm{d}\omega_1 \cdots \mathrm{d}\omega_d$ appearing in Equation (135). In terms of that (original) bound, we get the following theorem.

**Theorem 8.** *Assume that the spectral density is of the squared exponential form* (68)*. Then we have*

$$\begin{aligned}
&\int_{\frac{\pi\,\hat{m}}{2L_1}}^{\infty} \cdots \int_{\frac{\pi\,\hat{m}}{2L_d}}^{\infty} S(\omega_1,\ldots,\omega_d)\,\mathrm{d}\omega_1 \cdots \mathrm{d}\omega_d \\
&\leq D''\,\frac{\exp(-\gamma\,d\,m^{2/d})}{m} \leq \frac{D''}{m},
\end{aligned} \tag{69}$$

*for some constants $D'', \gamma > 0$ (which depend on $d$ and $L$).*

*Proof.* Due to separability of the spectral density, we have

$$\begin{aligned}
&\int_{\frac{\pi\,\hat{m}}{2L_1}}^{\infty} \cdots \int_{\frac{\pi\,\hat{m}}{2L_d}}^{\infty} S(\omega_1,\ldots,\omega_d)\,\mathrm{d}\omega_1 \cdots \mathrm{d}\omega_d \\
&= \left(s^2\,\sqrt{2\pi}\,\ell\right)^d \prod_{i=1}^{d} \int_{\frac{\pi\,\hat{m}}{2L_i}}^{\infty} \exp\left(-\frac{\ell^2\,\omega_i^2}{2}\right)\,\mathrm{d}\omega_i \\
&\leq \left(s^2\,\sqrt{2\pi}\,\ell\right)^d \left[\int_{\frac{\pi\,\hat{m}}{2L}}^{\infty} \exp\left(-\frac{\ell^2\,\omega_i^2}{2}\right)\,\mathrm{d}\omega_i\right]^d,
\end{aligned} \tag{70}$$

where $L = \min_k L_k$. By using the bound from Feller (1968), Section VII.1, Lemma 2, we get that is this

$$\begin{aligned}
&\leq \left(s^2\,\sqrt{2\pi}\,\ell^2\right)^d \left[\exp\left(-\frac{1}{2}\left[\frac{\pi\,\hat{m}}{2L\,\ell}\right]^2\right)\frac{2L\,\ell}{\pi\,\hat{m}}\right]^d \\
&= D''\,\frac{\exp(-\gamma\,d\,m^{2/d})}{m}.
\end{aligned} \tag{71}$$

$\square$

The above theorem tells that the convergence in the squared exponential case is faster than $\sim 1/m$, independently of the dimensionality $d$. It is worth noting though that the bound is not independent of the dimensionality in the sense that the constants do depend

on it. Strictly speaking, the convergence rate is $h(d)/m$, for some function $h$ which depends on $d$. However, as function of $m$, this rate is independent of the dimensionality.

## 4.4 Other Domains

It would also be possible carry out similar convergence analysis, for example, in a spherical domain. In that case the technical details become slightly more complicated, because instead sinusoidals we will have Bessel functions and the eigenvalues no longer form a uniform grid. This means that instead of Riemann integrals we need to consider weighted integrals where the distribution of the zeros of Bessel functions is explicitly accounted for. It might also be possible to use some more general theoretical results from mathematical analysis to obtain the convergence results. However, due to these technical challenges more general convergence proof will be developed elsewhere.

There is also a similar technical challenge in the analysis when the basis functions are formed by assuming an input density (see Section 2.4) instead of a bounded domain. Because explicit expressions for eigenfunctions and eigenvalues cannot be obtained in general, the elementary proof methods which we used here cannot be applied. Therefore the convergence analysis of this case is also left as a topic for future research.

## 5 Relationship to Other Methods

In this section we compare our method to existing sparse GP methods from a theoretical point of view. We consider two different classes of approaches: a class of inducing input methods based on the Nyström approximation (following the interpretation of Quiñonero-Candela and Rasmussen, 2005b; Bui et al, 2017), and direct spectral approximations.

## 5.1 Methods from the Nyström Family

A crude but rather effective scheme for approximating the eigendecomposition of the Gram matrix is the Nyström method (see, *e.g.*, Baker, 1977, for the integral approximation scheme). This method is based on choosing a set of $m$ inducing inputs $\mathbf{x}_\mathrm{u}$ and scaling the corresponding eigendecomposition of their corresponding covariance matrix $\mathbf{K}_\mathrm{u,u}$ to match that of the actual covariance. The Nyström approximations to the

$j$th eigenvalue and eigenfunction are

$$\widetilde{\lambda}_j = \frac{1}{m}\,\lambda_{\mathrm{u},j}, \tag{72}$$

$$\widetilde{\phi}_j(\mathbf{x}) = \frac{\sqrt{m}}{\lambda_{\mathrm{u},j}}\,k(\mathbf{x},\mathbf{x}_\mathrm{u})\,\phi_{\mathrm{u},j}, \tag{73}$$

where $\lambda_{\mathrm{u},j}$ and $\phi_{\mathrm{u},j}$ correspond to the $j$th eigenvalue and eigenvector of $\mathbf{K}_\mathrm{u,u}$. This scheme was originally introduced to the GP context by Williams and Seeger (2001). They presented a sparse scheme, where the resulting approximate prior covariance over the latent variables is $\mathbf{K}_\mathrm{f,u}\mathbf{K}_\mathrm{u,u}^{-1}\mathbf{K}_\mathrm{u,f}$, which can be derived directly from Equations (72) and (73).

As discussed by Quiñonero-Candela and Rasmussen (2005b), the Nyström method by Williams and Seeger (2001) does not correspond to a well-formed probabilistic model. However, several methods modifying the inducing point approach are widely used. The *Subset of Regressors* (SOR, Smola and Bartlett, 2001) method uses the Nyström approximation scheme for approximating the whole covariance function,

$$k_{\mathrm{SOR}}(\mathbf{x},\mathbf{x}') = \sum_{j=1}^{m} \widetilde{\lambda}_j\,\widetilde{\phi}_j(\mathbf{x})\,\widetilde{\phi}_j(\mathbf{x}'), \tag{74}$$

whereas the sparse Nyström method (Williams and Seeger, 2001) only replaces the training data covariance matrix. The SOR method is in this sense a complete Nyström approximation to the full GP problem. A method in-between is the *Deterministic Training Conditional* (DTC, Csató and Opper, 2002; Seeger et al, 2003) method which retains the true covariance for the training data, but uses the approximate cross-covariances between training and test data. For DTC, tampering with the covariance matrix causes the result not to actually be a Gaussian process. The *Variational Approximation* (VAR, Titsias, 2009) method modifies the DTC method by an additional trace term in the likelihood that comes from the variational bound.

The *Fully Independent (Training) Conditional* (FIC, Quiñonero-Candela and Rasmussen, 2005b) method (originally introduced as *Sparse Pseudo-Input GP* by Snelson and Ghahramani, 2006) is also based on the Nyström approximation but contains an additional diagonal term replacing the diagonal of the approximate covariance matrix with the values from the true covariance. The corresponding prior covariance function for FIC, is thus

$$\begin{aligned} &k_{\mathrm{FIC}}(\mathbf{x}_i,\mathbf{x}_j) \\ &= k_{\mathrm{SOR}}(\mathbf{x}_i,\mathbf{x}_j) + \delta_{i,j}(k(\mathbf{x}_i,\mathbf{x}_j) - k_{\mathrm{SOR}}(\mathbf{x}_i,\mathbf{x}_j)), \end{aligned} \tag{75}$$

where $\delta_{i,j}$ is the Kronecker delta.

Figure 3 illustrates the effect of the approximations compared to the exact correlation structure in the GP. The dashed contours show the exact correlation contours computed for three locations with the squared exponential covariance function. Figure 3a shows the results for the FIC approximation with 16 inducing points (locations shown in the figure). It is clear that the number of inducing points or their locations are not sufficient to capture the correlation structure. For similar figures and discussion on the effects of the inducing points, see Vanhatalo et al (2010). This behavior is not unique to SOR or FIC, but applies to all the methods from the Nyström family.

## 5.2 Direct Spectral Methods

The spectral analysis and series expansions of Gaussian processes has a long history. A classical result (see, *e.g.*, Loève, 1963; Van Trees, 1968; Adler, 1981; Cramér and Leadbetter, 2013, and references therein) is that in a compact set $\mathbf{x}, \mathbf{x}' \in \Omega \subset \mathbb{R}^d$ defined continuous covariance function can be expanded into a Mercer series

$$K(\mathbf{x}, \mathbf{x}') = \sum_j \gamma_j \, \varphi_j(\mathbf{x}) \, \varphi_j(\mathbf{x}'), \tag{76}$$

where $\gamma_j$ and $\varphi_j$ are the eigenvalues and the orthonormal eigenfunctions of the covariance function, respectively, defined as

$$\int_\Omega K(\mathbf{x}, \mathbf{x}') \, \varphi_j(\mathbf{x}') \, \mathrm{d}\mathbf{x}' = \gamma_j \, \varphi_j(\mathbf{x}). \tag{77}$$

Furthermore, the convergence happens absolutely and uniformly (Adler, 1981). This also means that we can approximate the covariance function with a finite truncation of the series and the approximation is guaranteed to converge to the exact covariance function when the number of terms is increased.

In the case of Gaussian processes we get that a zero mean Gaussian process with the covariance function $K(\mathbf{x}, \mathbf{x}')$ has the following Karhunen–Loeve series expansion in the domain $\Omega$:

$$f(\mathbf{x}) = \sum_j f_j \, \varphi_j(\mathbf{x}), \tag{78}$$

where $f_j$ are independent zero-mean Gaussian random variables with variances $\gamma_j$. The (also classical) generalization of this classical result to more general inner products was already discussed in Section 2.4.

In the case that $\Omega$ is not compact, but covers the whole $\mathbb{R}^d$, and when the covariance function is homogeneous, the eigenvalues defined by (77) are no longer discrete, but they can only be expressed as the spectral density $S(\omega)$ which can be seen as a continuum of

eigenvalues. The eigenfunctions become complex exponentials, that is, sines and cosines – which in turn are a subset of eigenfunctions of Laplace operator. In this background, what (20) essentially says is that we can approximate the Mercer expansion (76) by using the basis consisting of the Laplacian eigenfunctions $\varphi_j(\mathbf{x}) \approx \phi_j(\mathbf{x})$ and point-wise evaluations of the spectral density at the Laplacian eigenvalues $\gamma_j \approx S(\sqrt{-\lambda_j})$.

Another related classical connection is to the works in the relationship of spline interpolation and Gaussian process priors (Wahba, 1978; Kimeldorf and Wahba, 1970; Wahba, 1990). In particular, it is well-known (see, *e.g.*, Wahba, 1990) that spline smoothing can be seen as Gaussian process regression with a specific choice of covariance function. The relationship of the spline regularization with Laplace operators then leads to series expansion representations that are closely related to the approximations considered here.

In more recent machine learning context, the sparse spectrum GP (SSGP) method introduced by Lázaro-Gredilla et al (2010) uses the spectral representation of the covariance function for drawing random samples from the spectrum. These samples are used for representing the GP on a trigonometric basis

$$\phi(\mathbf{x}) = \big( \cos(2\pi \, \mathbf{s}_1^\mathsf{T} \mathbf{x}) \; \sin(2\pi \, \mathbf{s}_1^\mathsf{T} \mathbf{x}) \; \ldots$$
$$\cos(2\pi \, \mathbf{s}_h^\mathsf{T} \mathbf{x}) \; \sin(2\pi \, \mathbf{s}_h^\mathsf{T} \mathbf{x})\big), \quad (79)$$

where the spectral points $\mathbf{s}_r, r = 1, 2, \ldots, h$ ($2h = m$), are sampled from the spectral density of the original stationary covariance function (following the normalization convention used in the original paper). The covariance function corresponding to the SSGP scheme is now of the form

$$k_{\mathrm{SSGP}}(\mathbf{x}, \mathbf{x}') = \frac{2\sigma^2}{m} \phi(\mathbf{x}) \, \phi^\mathsf{T}(\mathbf{x}')$$
$$= \frac{\sigma^2}{h} \sum_{r=1}^h \cos\big(2\pi \, \mathbf{s}_r^\mathsf{T} (\mathbf{x} - \mathbf{x}')\big), \tag{80}$$

where $\sigma^2$ is the magnitude scale hyperparameter. This representation of the sparse spectrum method converges to the full GP in the limit of the number of spectral points going to infinity, and it is the preferred formulation of the method in one or two dimensions (see Lázaro-Gredilla, 2010, for discussion). We can interpret the SSGP method in (80) as a Monte Carlo approximation of the Wiener–Khintchin integral. In order to have a representative sample of the spectrum, the method typically requires the number of spectral points to be large. For high-dimensional inputs the number of required spectral points becomes overwhelming, and optimizing the spectral locations along with the hyperparameters attractive. However, as argued by Lázaro-Gredilla et al (2010), this option does not converge to
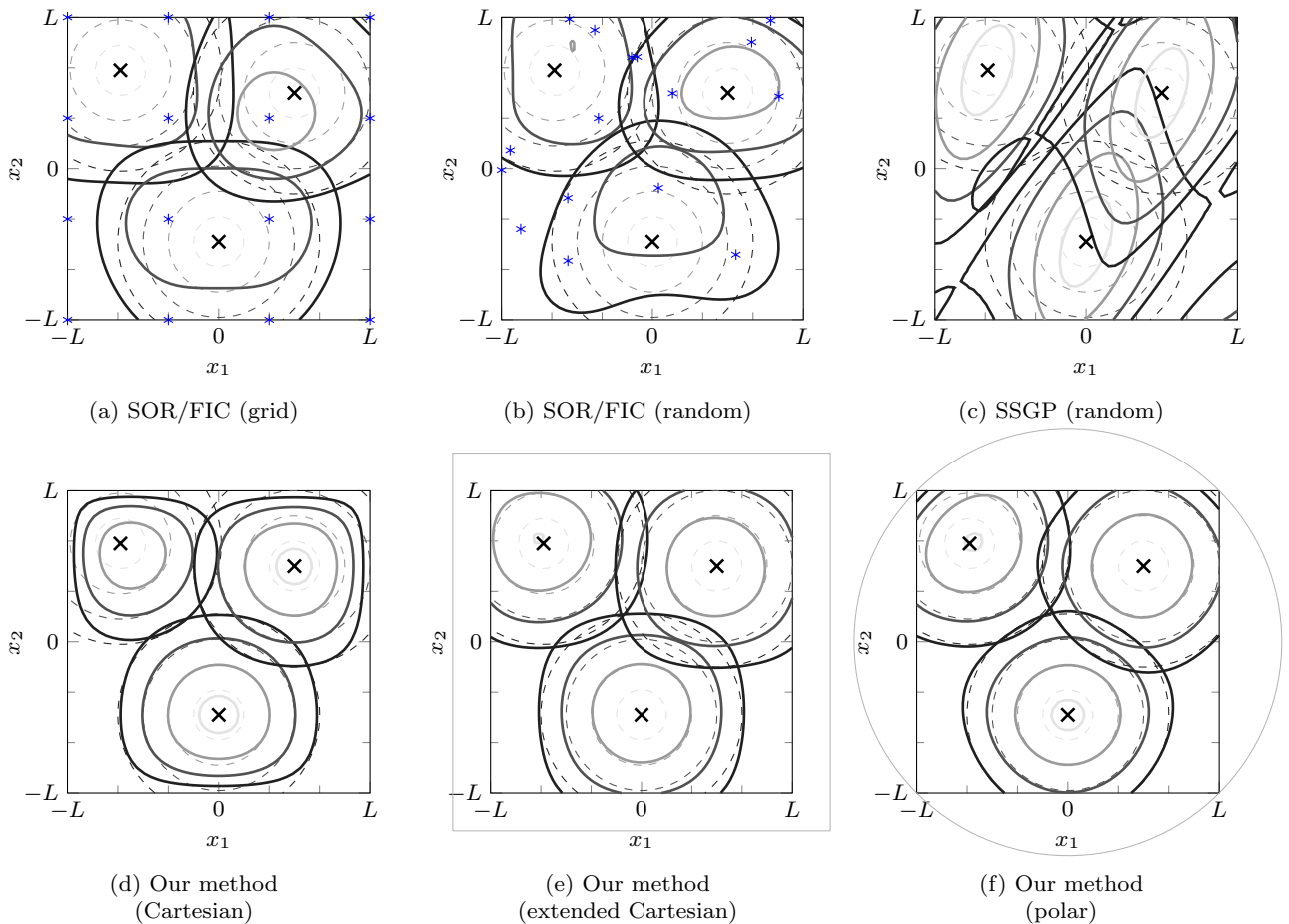
(a) SOR/FIC (grid)

(b) SOR/FIC (random)

(c) SSGP (random)

(d) Our method
(Cartesian)

(e) Our method
(extended Cartesian)

(f) Our method
(polar)

Fig. 3: Correlation contours computed for three locations (✗) corresponding to the squared exponential covariance function (exact contours dashed). The rank of each approximation is $m = 16$, and the locations of the inducing inputs are marked with blue stars (∗). The hyperparameters are the same in each figure. The domain boundary is shown in thin grey (——) if extended outside the box.

the full GP and suffers from overfitting to the training data (see Gal and Turner, 2015, for discussion on overfitting).

Contours for the sparse spectrum SSGP method are visualized in Figure 3c. Here the spectral points were chosen at random following Lázaro-Gredilla (2010). Because the basis functions are spanned using both sines and cosines, the number of spectral points was $h = 8$ in order to match the rank $m = 16$. These results agree well with those presented in the Lázaro-Gredilla et al (2010) for a one-dimensional example. For this particular set of spectral points some directions of the contours happen to match the true values very well, while other directions are completely off. Increasing the rank from 16 to 100 would give comparable results to the other methods.

Recently Hensman et al (2018) presented a variational Fourier feature approximation for Gaussian processes that was derived for the Matérn class of kernels,

where the approximation structure is set up by a low-rank plus diagonal structure. The key differences here are the fully diagonal (independent) structure in the $\mathbf{K}_{u,u}$ matrix (giving rise to additional speed-up) and the generality of only requiring the spectral density function to be known.

While SSGP is based on a sparse spectrum, the reduced-rank method proposed in this paper aims to make the spectrum as 'full' as possible at a given rank. While SSGP can be interpreted as a Monte Carlo integral approximation, the corresponding interpretation to the proposed method would be a numerical quadrature-based integral approximation (*cf.* the convergence proof in Appendix A.2). Figure 3d shows the same contours obtained by the proposed reduced-rank method. Here the eigendecomposition of the Laplace operator has been obtained for the square $\Omega = [-L, L] \times [-L, L]$ with Dirichlet boundary conditions. The contours match well with the full solution towards the middle of the domain.

The boundary effects drive the process to zero, which is seen as distortion near the edges.

Figure 3e shows how extending the boundaries just by 25% and keeping the number of basis functions fixed at 16, gives good results. The last Figure 3f corresponds to using a disk shaped domain instead of the rectangular. The eigendecomposition of the Laplace operator is done in polar coordinates, and the Dirichlet boundary is visualized by a circle in the figure.

## 5.3 Structure Exploiting and Decomposition Methods

Other methods for scalable Gaussian processes include many structure exploiting techniques that, similarly to general inducing input methods, aim to be agnostic to the choice of covariance function. They rather exploit the structure of the inputs (see Saatçi, 2012, for discussion on Kronecker and Toeplitz algebra), and not the GP prior per se. Most notably, Scalable Kernel Interpolation (SKI, Wilson and Nickisch, 2015) is an inducing point method that achieves $O(n + m \log m)$ time complexity and $\mathcal{O}(n + m)$ space complexity. Through local cubic kernel interpolation, the SKI framework is used in KISS-GP (see Wilson and Nickisch, 2015, for details) which uses Kronecker and Toeplitz algebra on grids of inducing inputs to speed up inference.

The computational complexity of the SKI approach scales cubically in the input dimenionality $d$. Other recent methods (*e.g.*, Gardner et al, 2018; Izmailov et al, 2018) have reduced the time complexity to linear in $d$ as well (*e.g.*, $\mathcal{O}(dn + dm \log m)$). These methods typically leverage parallelization (well suited for GPU calculations) or iterative methods.

Furthermore, general methods form numerical linear algebra for approximately solving eigenvalue and singular value problems allow for fast low-rank decompositions. These methods ignore the kernel learning perspective, but can provide useful tools in practice. For example, the pivoted Cholesky decomposition (Harbrecht et al, 2012; Bach, 2013) allows constructing a low-rank approximation to an $n \times n$ positive definite matrix in $O(nm^2)$ time. There are also methods for fast randomized singular value decompositions based on subsampled Hadamard transformations (*e.g.*, Boutsidis and Gittens, 2013), with some further details in Le et al (2013). These methods provide speedup to the general linear algebraic problem, but ignore the well structured nature of the specific application to Gaussian process regression with stationary prior covariance functions.
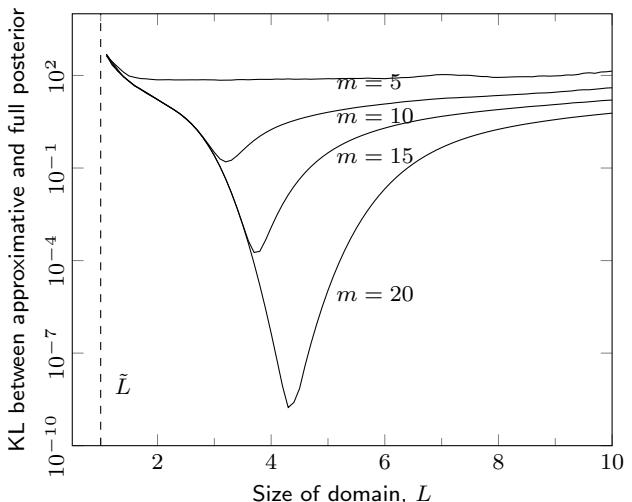


Fig. 4: The Kullback–Leibler divergence between the approximative and exact GP posterior by varying the boundary $L$ and keeping all other parameters fixed.

## 6 Experiments

In this section we aim to test the convergence results of the method in practice, provide examples of the practical use of the proposed method, and compare it against other methods that are typically used in a similar setting. We start with small simulated one-dimensional datasets, and then provide more extensive comparisons by using real-world data. We also consider an example of data, where the input domain is the surface of a sphere, and conclude our comparison by using a very large dataset to demonstrate what possibilities the computational benefits open.

## 6.1 Variation of Domain Size

In addition to the theoretical analysis of approximation error, we provide a study of the effect of choosing the domain size. We set up an experiment where we simulate data ($n = 100$ and all results averaged over 10 independent draws) from GP priors with a squared exponential covariance function with unit hyperparameters and corrupting additive Gaussian noise with variance $\sigma_\mathrm{n}^2 = 0.1^2$. The inputs are chosen uniformly randomly in $[-\tilde{L}, \tilde{L}]$ with $\tilde{L} = 1$. We study the effect of varying the boundary location $L \in (1, 10]$.

Figure 4 shows the Kullback–Leibler (KL) divergence (see, *e.g.*, Rasmussen and Williams, 2006, Appendix A for the the identities for the KL between two multivariate Gaussians) between the approximative GP posterior and the exact GP posterior evaluated over ten uniformly spaced points. The same curve is recalculated

(a) Gaussian process regression solution
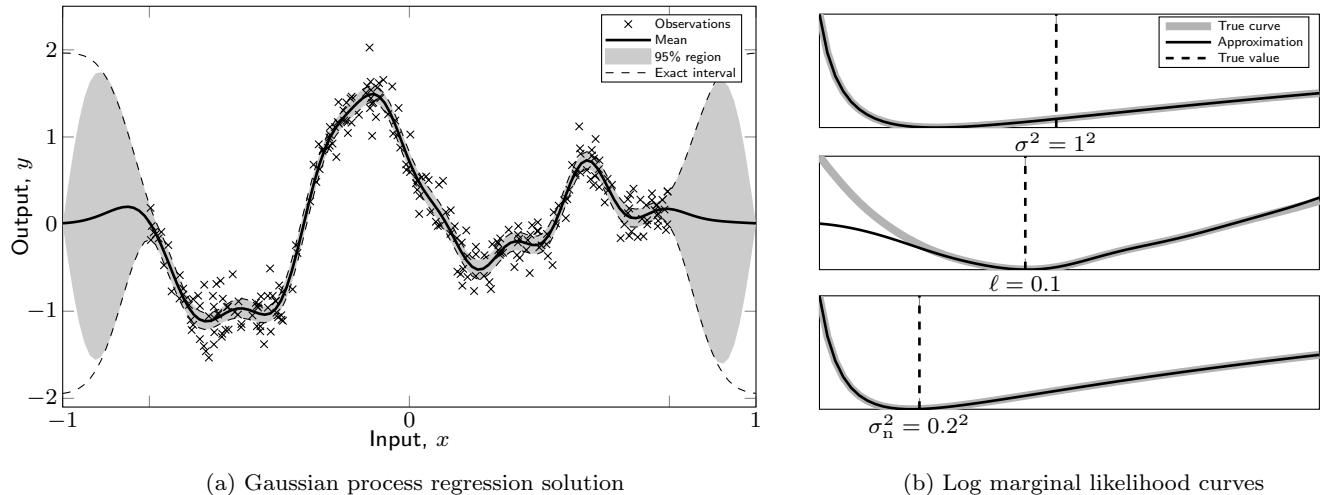


(b) Log marginal likelihood curves

Fig. 5: (a) 256 data points generated from a GP with hyperparameters $(\sigma^2, \ell, \sigma_n^2) = (1^2, 0.1, 0.2^2)$, the full GP solution, and an approximate solution with $m = 32$. (b) Negative marginal likelihood curves for the signal variance $\sigma^2$, length-scale $\ell$, and noise variance $\sigma_n^2$.

for $m = 5, 10, 15$, and 20. The figure shows that the KL has a single minimum that describes the trade-off of being far enough from the data but close enough not to start losing representative power with the given number of basis functions $m$. Even though the KL suggests there would be a single best choice for $L$, the practical sensitivity to the choice of $L$ is low. Already for $m = 5$, the MSE in the posterior mean is $10^{-5}$ (note that the data has unit magnitude scale) when $L$ is chosen one to two length-scales from the data boundary $\tilde{L}$.

### 6.2 Comparison Study

For assessing the performance of different methods we use 10-fold cross-validation and evaluate the following measures based on the validation set: the *standardized mean squared error* (SMSE) and the *mean standardized log loss* (MSLL), respectively defined as:

$$\text{SMSE} = \sum_{i=1}^{n_*} \frac{(y_{*i} - \mu_{*i})^2}{\text{Var}[y]}, \qquad (81)$$

and

$$\text{MSLL} = \frac{1}{2n_*} \sum_{i=1}^{n_*} \left( \frac{(y_{*i} - \mu_{*i})^2}{\sigma_{*i}^2} + \log 2\pi\sigma_{*i}^2 \right), \qquad (82)$$

where $\mu_{*i} = \mathbb{E}[f(\mathbf{x}_{*i})]$ and $\sigma_{*i}^2 = \mathbb{V}[f(\mathbf{x}_{*i})] + \sigma_n^2$ are the predictive mean and variance for test sample $i = 1, 2, \ldots, n_*$, and $y_{*i}$ is the actual test value. The training data variance is denoted by $\text{Var}[y]$. For all experiments, the values reported are averages over ten repetitions.

We compare our solution to SOR, DTC, VAR, and FIC using the implementations

provided in the GPstuff software package (version 4.3.1, see Vanhatalo et al, 2013) for Mathworks Matlab. The sparse spectrum SSGP method (Lázaro-Gredilla et al, 2010) was implemented into the GPstuff toolbox for the comparisons.[1] The reference implementation was modified such that also non-ARD covariances could be accounted for.

The $m$ inducing inputs for SOR, DTC, VAR, and FIC were chosen at random as a subset from the training data and kept fixed between the methods. For low-dimensional inputs, this tends to lead to good results and avoid over-fitting to the training data, while optimizing the input locations alongside hyperparameters becomes the preferred approach in high input dimensions (Quiñonero-Candela and Rasmussen, 2005b). The results are averaged over ten repetitions in order to present the average performance of the methods. In Sections 6.2 and 6.3, we used a Cartesian domain with Dirichlet boundary conditions for the new reduced-rank method. To avoid boundary effects, the domain was extended by 10% outside the inputs in each direction.

In the comparisons we followed the guidelines given by Chalupka et al (2013) for making comparisons between the actual performance of different methods. For hyperparameter optimization we used the `fminunc` routine in Matlab with a Quasi-Newton optimizer. We also tested several other algorithms, but the results were not sensitive to the choice of optimizer. The optimizer was run with a termination tolerance of $10^{-5}$ on the target function value and on the optimizer inputs. The num-

---

[1] The implementation is based on the code available from Miguel Lázaro-Gredilla: `http://www.tsc.uc3m.es/~miguel/downloads.php`.

(a) SMSE for the toy data



(b) MSLL for the toy data



(c) SMSE for the precipitation data



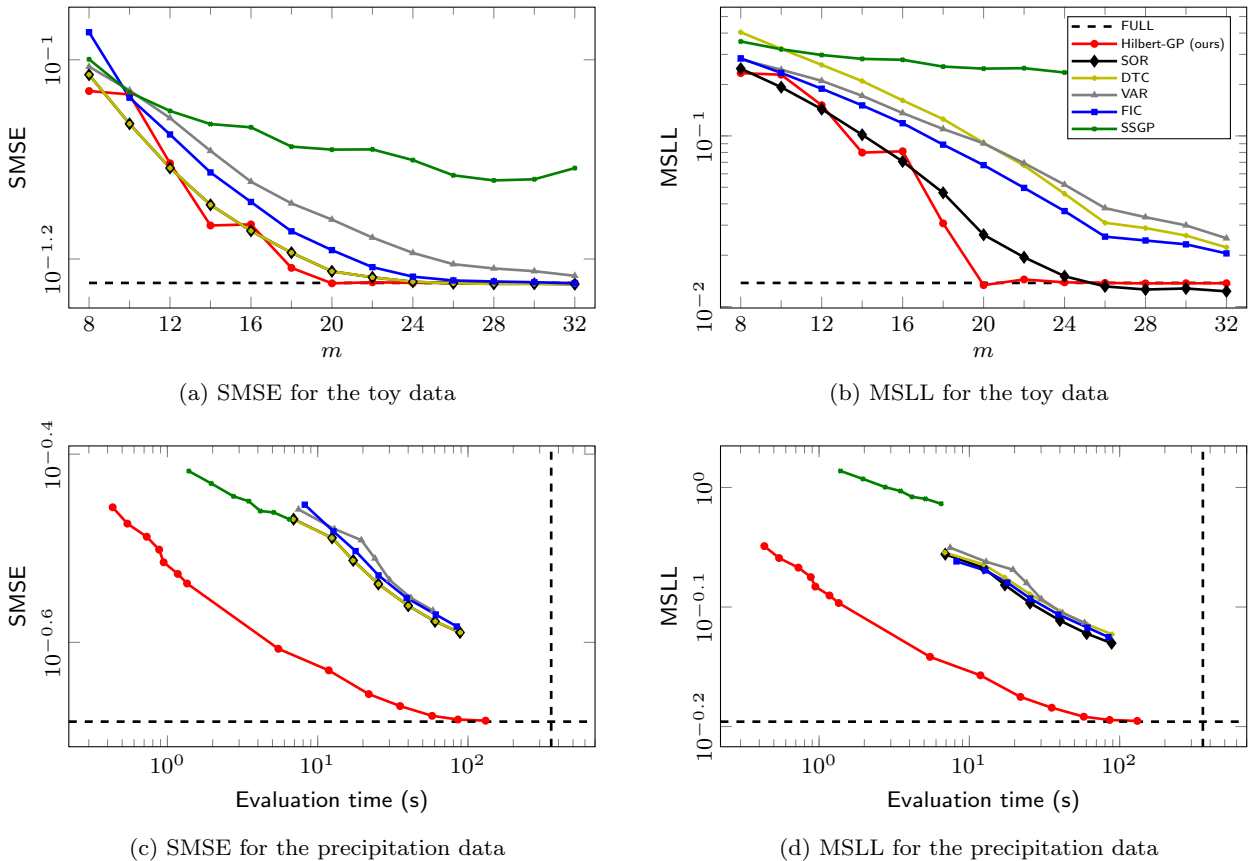(d) MSLL for the precipitation data

Fig. 6: Standardized mean squared error (SMSE) and mean standardized log loss (MSLL) results for the toy data ($d = 1$, $n = 256$) from Figure 5 and the precipitation data ($d = 2$, $n = 5776$) evaluated by 10-fold cross-validation and averaged over ten repetitions. The evaluation time includes hyperparameter learning.

ber of required target function evaluations stayed fairly constant for all the comparisons, making the comparisons for the hyperparameter learning bespoke.

Figure 5 shows a simulated example, where 256 data points are drawn from a Gaussian process prior with a squared exponential covariance function. We use the same parametrization as Rasmussen and Williams (2006) and denote the signal variance $\sigma^2$, length-scale $\ell$, and noise variance $\sigma_n^2$. Figure 5b shows the negative marginal log likelihood curves both for the full GP and the approximation with $m = 32$ basis functions. The likelihood curve approximations are almost exact and only differs from the full GP likelihood for small length-scales (roughly for values smaller than $2L/m$). Figure 5a shows the approximate GP solution. The mean estimate follows the exact GP mean, and the shaded region showing the 95% confidence area differs from the exact solution (dashed) only near the boundaries.

Figures 6a and 6b show the SMSE and MSLL values for $m = 8, 10, \ldots, 32$ inducing inputs and basis functions for the toy dataset from Figure 5. The convergence of the proposed reduced rank method is fast and

as soon as the number of eigenfunctions is large enough ($m = 20$) to account for the short length-scales, the approximation converges to the exact full GP solution (shown by the dashed line).

In this case the SOR method that uses the Nyström approximation to directly approximate the spectrum of the full GP (see Section 5) seems to give good results. However, as the resulting approximation in SOR corresponds to a singular Gaussian distribution, the predictive variance is underestimated. This can be seen in Figure 6b, where SOR seems to give better results than the full GP. These results are however due to the smaller predictive variance on the test set. DTC tries to fix this shortcoming of SOR—they are identical in other respects except predictive variance evaluation—and while SOR and DTC give identical results in terms of SMSE, they differ in MSLL. We also note that additional trace term in the marginal likelihood in VAR makes the likelihood surface flat, which explains the differences in the results in comparison to DTC.

The sparse spectrum SSGP method did not perform well on average. Still, it can be seen that it converges

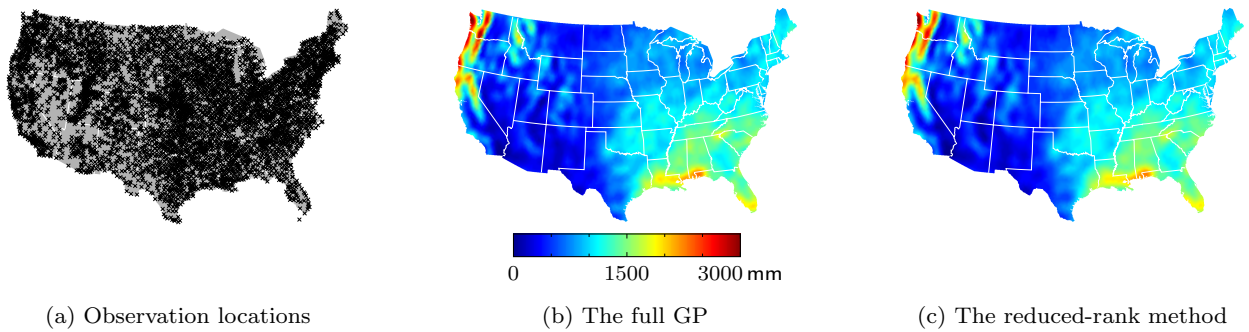(a) Observation locations  (b) The full GP  (c) The reduced-rank method

Fig. 7: Interpolation of the yearly precipitation levels using reduced-rank GP regression. Subfigure 7a shows the $n = 5776$ weather station locations. Subfigures 7b and 7c show the results for the full GP model and the new reduced-rank GP method.

towards the performance of the full GP. The dependence on the number of spectral points differs from the rest of the methods, and a rank of $m = 32$ is not enough to meet the other methods. However, in terms of best case performance over the ten repetitions with different inducing inputs and spectral points, both FIC and SSGP outperformed SOR, DTC, and VAR. Because of its 'dense spectrum' approach, the proposed reduced-rank method is not sensitive to the choice of spectral points, and thus the performance remained the same between repetitions. In terms of variance over the 10-fold cross-validation folds, the methods in order of growing variance in the figure legend (the variance approximately doubling between FULL and SSGP).

### 6.3 Precipitation Data

As a real-data example, we consider a precipitation data set that contain US annual precipitation summaries for year 1995 ($d = 2$ and $n = 5776$, available online, see Vanhatalo et al, 2013). The observation locations are shown on a map in Figure 7a.

We limit the number of inducing inputs and spectral points to $m = 128, 192, \ldots, 512$. For the our Hilbert-GP method we additionally consider ranks $m = 1024, 1536, \ldots, 4096$, and show that this causes a computational burden of the same order as the conventional sparse GP methods with smaller $m$s. To avoid boundary effects, the domain was extended by 10% outside the inputs in each direction.

In order to demonstrate the computational benefits of the proposed model, we also present the running time of the GP inference (including hyperparameter optimization). All methods were implemented under a similar framework in the GPstuff package, and they all employ similar reformulations for numerical stability. The key difference in the evaluation times comes from

hyperparameter optimization, where SOR, DTC, VAR, FIC, and SSGP scale as $\mathcal{O}(nm^2)$ for each evaluation of the marginal likelihood. The proposed reduced-rank method scales as $\mathcal{O}(m^3)$ for each evaluation (after an initial cost of $\mathcal{O}(nm^2)$).

Figures 6c and 6d show the SMSE and MSLL results for this data against evaluation time. On this scale we note that the evaluation time and accuracy, both in terms of SMSE and MSLL, are alike for SOR, DTC, VAR, and FIC. SSGP is faster to evaluate in comparison with the Nyström family of methods, which comes from the simpler structure of the approximation. Still, the number of required spectral points to meet a certain average error level is larger for SSGP.

The results for the proposed reduced-rank method (Hilbert-GP) show that with two input dimensions, the required number of basis functions is larger. For the first seven points, we notice that even though the evaluation is two orders of magnitude faster, the method performs only slightly worse in comparison to conventional sparse methods. By considering higher ranks (the next seven points), our method converges to the performance of the full GP (both in SMSE and MSLL), while retaining a computational time comparable to the conventional methods. This type of spatial medium-size GP regression problems can thus be solved in seconds.

Figures 7b and 7c show interpolation of the precipitation levels using a full GP model and the reduced-rank method ($m = 1728$), respectively. The results are practically identical, as is easy to confirm from the color surfaces. Obtaining the reduced-rank result (including initialization and hyperparameter learning) took slightly less than 30 seconds on a laptop computer (MacBook Air, Late 2010 model, 2.13 GHz, 4 GB RAM), while the full GP inference took approximately 18 minutes.

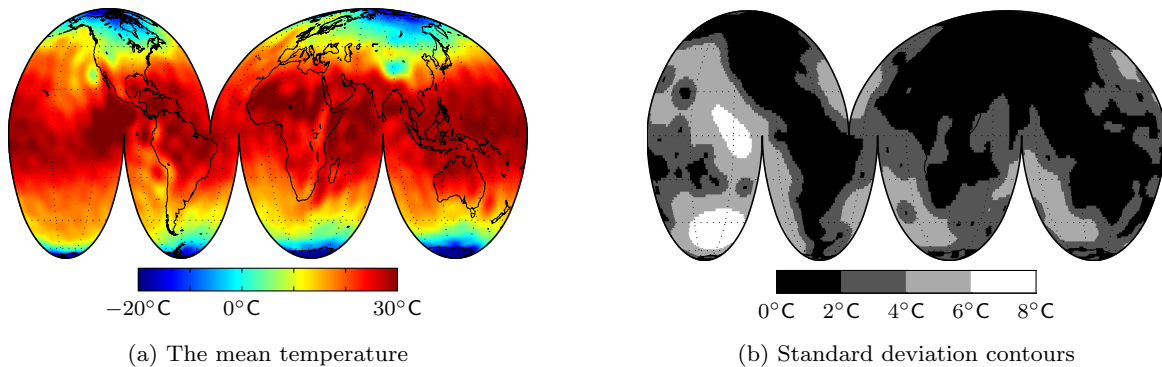(a) The mean temperature



(b) Standard deviation contours

Fig. 8: Modeling of the yearly mean temperature on the spherical surface of the Earth ($n = 11\,028$). Figure 8b shows the standard deviation contours which match well with the continents.

## 6.4 Temperature Data on the Surface of the Globe

We also demonstrate the use of the method in non-Cartesian coordinates. We consider modeling of the spatial mean temperature over a number of $n = 11\,028$ locations around the globe.[2]

As earlier demonstrated in Figure 2, we use the Laplace operator in spherical coordinates as defined in (31). The eigenfunctions for the angular part are the Laplace's spherical harmonics. The evaluation of the approximation does not depend on the coordinate system, and thus all the equations presented in the earlier sections remain valid. We use the squared exponential covariance function and $m = 1089$ basis functions.

Figure 8 visualizes the modeling outcome. The results are visualized using an interrupted projection (an adaption of the Goode homolosine projection) in order to preserve the length-scale structure across the map. The uncertainty is visualized in Figure 8b, which corresponds to the $n = 11,028$ observation locations that are mostly spread over the continents and western countries (the white areas in Figure 8b contain no observations). Obtaining the reduced-rank result (including initialization and hyperparameter learning) took approximately 50 seconds on a laptop computer (MacBook Air, Late 2010 model, 2.13 GHz, 4 GB RAM), which scales with $n$ in comparison to the evaluation time in the previous section.

## 6.5 Additive Modelling of Airline Delays

In order to fully use the computational benefits and also underline a way of applying the method to high-dimensional inputs, we consider a large dataset for pre-dicting airline delays. The US flight delay prediction example (originally considered by Hensman et al, 2013) is a standard test data set in Gaussian process regression. This is due to the clearly non-stationary behavior and its massive size, with nearly 6 million records.

We aim to replicate and extend to the results previously presented in the work by Hensman et al (2018) for the Variational Fourier Features (VFF) method. This example has also been used by Deisenroth and Ng (2015), where it was solved using distributed Gaussian processes, and by Samo and Roberts (2016) who use this example for demonstrating the computational efficiency of string Gaussian processes. Adam et al (2016) used this data set as an example where the model can be formed by the addition of multiple underlying components.

The data consists of flight arrival and departure times for every commercial flight in the USA for the year 2008. We use the standard eight covariates $\mathbf{x}$ (see Hensman et al, 2013) which are the age of the aircraft (number of years since deployment), route distance, airtime, departure time, arrival time, day of the week, day of the month, and month. The target is to predict the delay of the aircraft at landing (in minutes), $y$.

This regression task is set up similarly as in Hensman et al (2018) and Adam et al (2016), as a Gaussian process regression model with a prior covariance structure given as a sum of ovariance functions for each input dimension and assuming the observations are corrupted by independent Gaussian noise, $\varepsilon_i \sim \mathcal{N}(0, \sigma_{\mathrm{n}}^2)$. The model is

$$f(\mathbf{x}) \sim \mathcal{GP}\left(0, \sum_{d=1}^{8} k_{\mathrm{se}}(x_d, x_d')\right),$$
$$y_i = f(\mathbf{x}_i) + \varepsilon_i, \tag{83}$$

---

[2] The data are available for download from US National Climatic Data Center: `http://www7.ncdc.noaa.gov/CDO/cdoselect.cmd` (accessed January 3, 2014).

Table 1: Predictive mean squared errors (MSEs) and negative log predictive densities (NLPDs) with one standard deviation on the airline arrival delays experiment (input dimensionality $d = 8$) for a number of data points ranging up to almost 6 million. The Hilbert-GP method is on par with the VFF method albeit being clearly faster due to the diagonalizable structure (solving the regression problem including hyperparameter optimization in 41 seconds on a CPU-only laptop computer).

| $n$ | 10,000 | | 100,000 | | 1,000,000 | | 5,929,413 | |
|---|---|---|---|---|---|---|---|---|
| | MSE | NLPD | MSE | NLPD | MSE | NLPD | MSE | NLPD |
| Hilbert-GP | $0.97 \pm 0.14$ | $1.404 \pm 0.071$ | $0.80 \pm 0.06$ | $1.311 \pm 0.038$ | $0.83 \pm 0.02$ | $1.329 \pm 0.011$ | $0.827 \pm 0.005$ | $1.324 \pm 0.003$ |
| VFF | $0.89 \pm 0.15$ | $1.362 \pm 0.091$ | $0.82 \pm 0.05$ | $1.319 \pm 0.030$ | $0.83 \pm 0.01$ | $1.326 \pm 0.008$ | $0.827 \pm 0.004$ | $1.324 \pm 0.003$ |
| SVIGP | $0.89 \pm 0.16$ | $1.354 \pm 0.096$ | $0.79 \pm 0.05$ | $1.299 \pm 0.033$ | $0.79 \pm 0.01$ | $1.301 \pm 0.009$ | $0.791 \pm 0.005$ | $1.300 \pm 0.003$ |
| Full-RBF | $0.89 \pm 0.16$ | $1.349 \pm 0.098$ | N/A | N/A | N/A | N/A | N/A | N/A |
| Full-additive | $0.89 \pm 0.16$ | $1.362 \pm 0.096$ | N/A | N/A | N/A | N/A | N/A | N/A |

for $i = 1, 2, \ldots, m$. We used $m = 40$ basis functions per input dimension. The boundary is set to a distance of two times the range of the data for each dimension.

We consider several subset sizes of the data, each selected uniformly at random: $n = 10,000$, $100,000$, $1,000,000$, and $5,929,413$ (all data). In each case, two thirds of the data is used for training and one third for testing. For each subset size the training is repeated ten times. The random splits are exactly the same as in Hensman et al (2018).

Table 1 shows the (normalized) predictive mean squared errors (MSEs) and the negative log predictive densities (NLPDs) with one standard deviation on the airline arrival delays experiment. The table shows that the Hilbert-GP method is directly on par with the Variational Fourier Features (VFF) method. For the smaller subsets some variability in the results is visible, even though the MSEs and NLPDs are within one standard deviation of one another for VFF and Hilbert-GP. For the data sets in the millions, VFF and Hilbert-GP perform practically equally well. Further analysis and interpretation of the data and model can be found in Hensman et al (2018). We have omitted reporting results for the String GP method Samo and Roberts (2016), the Bayesian committee machine (BCM, Tresp, 2000), and the robust Bayesian committee machine (rBCM, Deisenroth and Ng, 2015). Each of these performed worse than any of the included methods, and the resulting numbers can be found listed in Hensman et al (2018) and Samo and Roberts (2016).

Running the Hilbert-GP method in this experiment (including hyperparameter training and prediction) with all 5.93 million data took $41 \pm 2$ seconds ($120 \pm 7$ s CPU time) on a MacBook Pro laptop (with all calculation done on the CPU). The is clearly faster than the VFF method with $265 \pm 6$ seconds ($626 \pm 11$ s CPU time), where our computational gain comes from the fully diagonal structure of the covariance. For comparison, the SVIGP method (Hensman et al, 2013) required $5.1 \pm 0.1$ hours of computing ($27.0 \pm 0.8$ h CPU

time) on a cluster. Samo and Roberts (2016) report that running the String GP took 91.0 hours total CPU time (or 15 h of wall-clock time on an 8-core machine). Izmailov et al (2018) also report results for the airline dataset, where one pass over the data taking 5200 seconds, when running on a Nvidia Tesla K80 GPU and not assuming additive structure.

### 6.6 Gaussian Process Driven Poisson Equation

As discussed in Section 3.4, our framework also directly extends to inverse problems and latent force models. As this final experiment, we demonstrate the use of the approximation in the latent force model (LFM)

$$-\nabla^2 g(\mathbf{x}) = f(\mathbf{x}),$$
$$y_i = g(\mathbf{x}_i) + \varepsilon_i, \tag{84}$$

where $\mathbf{x} \in \mathbb{R}^2$ and $f(\mathbf{x}) \sim \mathcal{GP}(0, k(\mathbf{x}, \mathbf{x}'))$ is the input with a squared exponential covariance function prior. This problem can also be interpreted as a inverse problem where the measurement operator is the Green's operator $\mathcal{H} = (-\nabla^2)^{-1}$:

$$y_i = (\mathcal{H}f)(\mathbf{x}_i) + \varepsilon_i. \tag{85}$$

If we assume that the boundary conditions of the problem are the same as we used for forming the basis functions in (10), then if we put $g(\mathbf{x}) \approx \sum_{j=1}^{m} g_j \, \phi_j(\mathbf{x})$, we get

$$-\nabla^2 g(\mathbf{x}) \approx -\sum_{j=1}^{m} g_j \, \nabla^2 \phi_j(\mathbf{x}) = \sum_{j=1}^{m} g_j \, \lambda_j \, \phi_j(\mathbf{x}) \tag{86}$$

and thus by further putting $f(\mathbf{x}) \approx \sum_{j=1}^{m} f_j \, \phi_j(\mathbf{x})$, the approximation to the equation $-\nabla^2 g(\mathbf{x}) = f(\mathbf{x})$ becomes

$$\sum_{j=1}^{m} g_j \, \lambda_j \, \phi_j(\mathbf{x}) = \sum_{j=1}^{m} f_j \, \phi_j(\mathbf{x}) \tag{87}$$
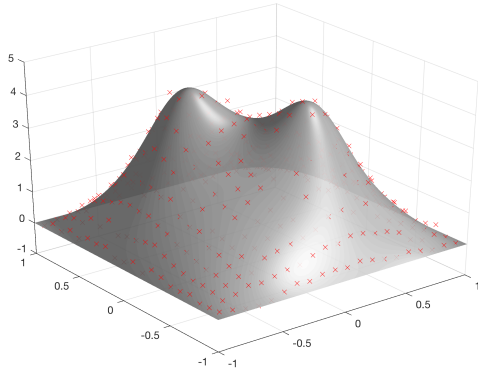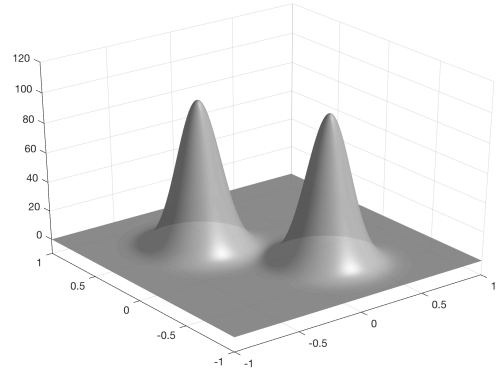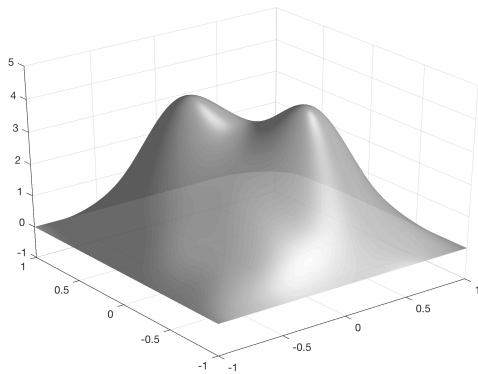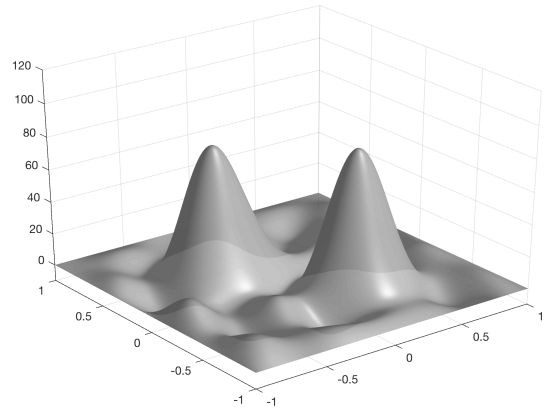
(a) The true solution $g(x_1, x_2)$ and the measurements.



(b) The true input $f(x_1, x_2)$.



(c) The estimate of solution $g(x_1, x_2)$.



(d) The estimate of input $f(x_1, x_2)$.

Fig. 9: Gaussian process inference on the Poisson equation.

which allows us to solve $f_j = g_j/\lambda_j$. This implies that we approximately have $(-\nabla^2)^{-1}\phi_j = \phi_j/\lambda_j$ which reduces Equations (52) to

$$(\mathcal{H}\,\mathcal{H}'\,k)(\mathbf{x}, \mathbf{x}') \approx \sum_j \lambda_j^{-2}\,S(\sqrt{\lambda_j})\,\phi_j(\mathbf{x})\,\phi_j(\mathbf{x}'),$$
$$(\mathcal{H}'k(\mathbf{x}_*, \cdot))(\mathbf{x}') \approx \sum_j \lambda_j^{-1}\,S(\sqrt{\lambda_j})\,\phi_j(\mathbf{x}_*)\,\phi_j(\mathbf{x}'), \tag{88}$$

after which we can proceed with (51). Alternatively we can directly use (53) with $\tilde{\mathbf{\Phi}}_{ij} = \phi_j(\mathbf{x}_i)/\lambda_j$.

Figure 9 shows the result of applying the proposed method to this model with the input function shown in Figure 9b. The true solution and the simulated measurements (with standard deviation of $1/10$) are shown in Figure 9a. The scale $\sigma^2$ and length scale $\ell$ of the SE covariance function were estimated by maximum likelihood method and the number of basis functions used for solving the GP regression problem was 100 (for simulation we used 255 basis functions). The estimates of

the input and the solution function are shown in Figures 9b and 9a, respectively. As can be seen in the figures, the estimate of the solution is very good, as can be expected from the fact that we obtain direct (although noisy) measurements from it. The estimate of the input is less accurate, but still approximates the true input well.

## 7 Conclusion and Discussion

In this paper we have proposed a novel approximation scheme for forming approximate eigendecompositions of covariance functions in terms of the Laplace operator eigenbasis and the spectral density of the covariance function. The eigenfunction decomposition of the Laplacian can easily be formed in various domains, and the eigenfunctions are independent of the choice of hyperparameters of the covariance.

An advantage of the method is that it has the ability to approximate the eigendecomposition using only the eigendecomposition of the Laplacian and the spectral density of the covariance function, both of which are closed-from expressions. This together with having the eigenvectors in $\boldsymbol{\Phi}$ mutually orthogonal and independent of the hyperparameters, is the key to efficiency. This allows an implementation with a computational cost of $\mathcal{O}(nm^2)$ (initial) and $\mathcal{O}(m^3)$ (marginal likelihood evaluation), with negligible memory requirements.

Of the infinite number of possible basis functions only an extremely small subset are of any relevance to the GP being approximated. In GP regression the model functions are conditioned on a covariance function (kernel), which imposes desired properties on the solutions. We choose the basis functions such that they are as close as possible (w.r.t. the Frobenius norm) to those of the particular covariance function. Our method gives the exact eigendecomposition of a GP that has been constrained to be zero at the boundary of the domain.

The method allows for theoretical analysis of the error induced by the truncation of the series and the boundary effects. This is something new in this context and extremely important, for example, in medical imaging applications. The approximative eigendecomposition also opens a range of interesting possibilities for further analysis. In *learning curve* estimation, the eigenvalues of the Gaussian process can now be directly approximated. For example, we can approximate the Opper–Vivarelli bound (Opper and Vivarelli, 1999) as

$$\epsilon_{\mathrm{OV}}(n) \approx \sigma_{\mathrm{n}}^2 \sum_j \frac{S(\sqrt{\lambda_j})}{\sigma_{\mathrm{n}}^2 + n\,S(\sqrt{\lambda_j})}. \tag{89}$$

Sollich's eigenvalue based bounds (Sollich and Halees, 2002) can be approximated and analyzed in an analogous way.

However, some of these abilities come with a cost. As demonstrated throughout the paper, restraining the domain to boundary conditions introduces edge effects. These are, however, known and can be accounted for. Extrapolating with a stationary covariance function outside the training inputs only causes the predictions to revert to the prior mean and variance. Therefore we consider the boundary effects a minor problem for practical use.

Although at first sight the method appears to have a bad (exponential) scaling with respect to the input dimensionality, as shown by the analysis in Section 4.3, this is not true. By increasing the differentiability order of the covariance function appropriately we can keep the convergence rate at the level $\sim 1/m^a$, for a given constant $a > 0$ and with total of $m$ terms in the

series, regardless of the input dimensionality. Furthermore, Theorem 8 shows that for squared exponential covariance function the convergence rate is always better than $\sim 1/m$, independently of the input dimensionality.

Further resources related to the proposed method and implementation details in form of code are available at `https://github.com/AaltoML/hilbert-gp`.

## A Proofs of Convergence Theorems

### A.1 Auxiliary Lemmas

In this section we present a few lemmas that will be needed in the proofs in the next sections. The lemmas are quite classical results on the convergence of Riemannian sums, but as it is hard to find exactly the same results in other literature, for completeness we prove the lemmas here.

**Lemma 9.** *Let* $\Delta > 0$ *and* $\alpha \in [0,1)$ *be given constants,* $m = 0, 1, 2, \ldots$ *some nonnegative integer, and assume that the* $f(\omega)$ *is a bounded integrable function defined on* $\omega \geq m\,\Delta$ *with bounded derivative on* $\omega > m\,\Delta$ *such that* $\int_{m\,\Delta}^{\infty} |f'(\omega)|\,d\omega = C^{(m)} < \infty$. *Then we have*

$$\left| \int_{m\,\Delta}^{\infty} f(\omega)\,d\omega - \sum_{j=m+1}^{\infty} f(j\,\Delta - \alpha\,\Delta)\,\Delta \right| \leq C^{(m)}\,\Delta. \tag{90}$$

*Furthermore, provided that* $\int_0^{\infty} |f'(\omega)|\,d\omega = C^{(0)} < \infty$, *this bound can be made independent of* $m$:

$$\left| \int_{m\,\Delta}^{\infty} f(\omega)\,d\omega - \sum_{j=m+1}^{\infty} f(j\,\Delta - \alpha\,\Delta)\,\Delta \right| \leq C^{(0)}\,\Delta. \tag{91}$$

*Proof.* We can write

$$\int_{m\,\Delta}^{\infty} f(\omega)\,d\omega = \sum_{j=m+1}^{\infty} \int_{(j-1)\,\Delta}^{j\,\Delta} f(\omega)\,d\omega. \tag{92}$$

By the fundamental theorem of calculus we get

$$f(\omega) = f(j\,\Delta - \alpha\,\Delta) + \int_{j\,\Delta - \alpha\,\Delta}^{\omega} f'(\omega)\,d\omega, \tag{93}$$

which gives for $\omega \in ((j-1)\,\Delta, j\,\Delta]$

$$
\begin{aligned}
|f(\omega) - f(j\,\Delta - \alpha\,\Delta)| &\leq \left| \int_{j\,\Delta - \alpha\,\Delta}^{\omega} f'(\omega)\, d\omega \right| \\
&\leq \left| \int_{j\,\Delta - \alpha\,\Delta}^{\omega} |f'(\omega)|\, d\omega \right| \\
&\leq \int_{(j-1)\,\Delta}^{j\,\Delta} |f'(\omega)|\, d\omega
\end{aligned}
\tag{94}
$$

We now get

$$
\begin{aligned}
& \left| \int_{m\,\Delta}^{\infty} f(\omega)\, d\omega - \sum_{j=m+1}^{\infty} f(j\,\Delta - \alpha\,\Delta)\,\Delta \right| \\
&= \left| \sum_{j=m+1}^{\infty} \int_{(j-1)\,\Delta}^{j\,\Delta} [f(\omega) - f(j\,\Delta - \alpha\,\Delta)]\, d\omega \right| \\
&\leq \sum_{j=m+1}^{\infty} \int_{(j-1)\,\Delta}^{j\,\Delta} |f(\omega) - f(j\,\Delta - \alpha\,\Delta)|\, d\omega \\
&\leq \sum_{j=m+1}^{\infty} \int_{(j-1)\,\Delta}^{j\,\Delta} \left[ \int_{(j-1)\,\Delta}^{j\,\Delta} |f'(\omega)|\, d\omega \right] d\omega \\
&= \sum_{j=m+1}^{\infty} \int_{(j-1)\,\Delta}^{j\,\Delta} |f'(\omega)|\, d\omega\, \Delta \\
&= \underbrace{\int_{m\,\Delta}^{\infty} |f'(\omega)|\, d\omega}_{C^{(m)}}\, \Delta \\
&\leq \underbrace{\int_{0}^{\infty} |f'(\omega)|\, d\omega}_{C^{(0)}}\, \Delta,
\end{aligned}
\tag{95}
$$

which concludes the proof.                                    □

**Lemma 10.** *Assume that $f(\omega)$ is a bounded integrable function defined on $\omega \geq 0$ with bounded derivative on $\omega > 0$ and $g(\omega)$ is a bounded function defined on $\omega \geq 0$ such that $|g(w)| \leq D$. Further assume that $\int_0^\infty |f'(\omega)|\, d\omega = C < \infty$. Then for any $\alpha, \beta \in [0,1)$ and $\Delta > 0$ we have*

$$
\left| \sum_{j=1}^{\infty} [f(j\,\Delta) - f(j\,\Delta - \alpha\,\Delta)]\, g(j\,\Delta - \beta\,\Delta) \right| \leq C\, D.
\tag{96}
$$

*Proof.* By using (94) with $\omega = j\,\Delta - \alpha\,\Delta$ we get

$$
|f(j\,\Delta) - f(j\,\Delta - \alpha\,\Delta)| \leq \int_{(j-1)\,\Delta}^{j\,\Delta} |f'(\omega)|\, d\omega,
\tag{97}
$$

and further

$$
\begin{aligned}
& \left| \sum_{j=1}^{\infty} [f(j\,\Delta) - f(j\,\Delta - \alpha\,\Delta)]\, g(j\,\Delta - \beta\,\Delta) \right| \\
&\leq \sum_{j=1}^{\infty} |f(j\,\Delta) - f(j\,\Delta - \alpha\,\Delta)|\, |g(j\,\Delta - \beta\,\Delta)| \\
&\leq \sum_{j=1}^{\infty} \int_{(j-1)\,\Delta}^{j\,\Delta} |f'(\omega)|\, d\omega\, D \\
&= \int_{0}^{\infty} |f'(\omega)|\, d\omega\, D \\
&= C\, D.
\end{aligned}
\tag{98}
$$

□

**Lemma 11.** *Assume that $f(\omega) \geq 0$ is a positive bounded integrable function defined on $\omega \geq 0$ with bounded derivative on $\omega > 0$ such that $\int_0^\infty f(\omega)\, d\omega = C_0 < \infty$ and $\int_0^\infty |f'(\omega)|\, d\omega = C_1 \leq \infty$, and $g(\omega)$ is a bounded integrable function defined on $\omega \geq 0$ with bounded derivative on $\omega > 0$ such that $|g'(\omega)| \leq D$. Then for any $\alpha, \beta \in [0,1)$ and $\Delta > 0$ we have for $C = C_1 + C_0$:*

$$
\left| \sum_{j=1}^{\infty} f(j\,\Delta - \alpha\,\Delta) [g(j\,\Delta) - g(j\,\Delta - \beta\,\Delta)] \right| \leq C\, D.
\tag{99}
$$

*Proof.* By applying the mean value theorem to (97) we get that for some $\omega_j^* \in [j\,\Delta - \alpha\,\Delta, j\,\Delta]$ we have

$$
|g(j\,\Delta) - g(j\,\Delta - \beta\,\Delta)| \leq |g'(\omega_j^*)|\,\beta\,\Delta \leq |g'(\omega_j^*)|\,\Delta \leq D\,\Delta.
\tag{100}
$$

By using Lemma 9 we get

$$
\begin{aligned}
& \sum_{j=1}^{\infty} f(j\,\Delta - \alpha\,\Delta)\,\Delta \\
&= \left| \sum_{j=1}^{\infty} f(j\,\Delta - \alpha\,\Delta)\,\Delta - \int_0^\infty f(\omega)\, d\omega + \int_0^\infty f(\omega)\, d\omega \right| \\
&\leq \left| \sum_{j=1}^{\infty} f(j\,\Delta - \alpha\,\Delta)\,\Delta - \int_0^\infty f(\omega)\, d\omega \right| + \left| \int_0^\infty f(\omega)\, d\omega \right| \\
&\leq C_1 + C_0 = C.
\end{aligned}
\tag{101}
$$

Hence,

$$
\begin{aligned}
& \left| \sum_{j=1}^{\infty} f(j\,\Delta - \alpha\,\Delta) [g(j\,\Delta) - g(j\,\Delta - \beta\,\Delta)] \right| \\
&\leq \sum_{j=1}^{\infty} f(j\,\Delta - \alpha\,\Delta)\, |g(j\,\Delta) - g(j\,\Delta - \beta\,\Delta)| \\
&\leq \sum_{j=1}^{\infty} f(j\,\Delta - \alpha\,\Delta)\,\Delta\, D \\
&\leq C\, D.
\end{aligned}
\tag{102}
$$

□

### A.2 Proof of Theorem 1

The Wiener–Khinchin identity and the symmetry of the spectral density allows us to write

$$
\begin{aligned}
k(x,x') &= \frac{1}{2\pi} \int_{-\infty}^{\infty} S(\omega)\, \exp(-i\,\omega\,(x-x'))\, d\omega \\
&= \frac{1}{\pi} \int_{0}^{\infty} S(\omega)\, \cos(\omega\,(x-x'))\, d\omega.
\end{aligned}
\tag{103}
$$

In a one-dimensional domain $\Omega = [-L, L]$ with Dirichet boundary conditions we have an $m$-term approximation of the form

$$
\begin{aligned}
& \widetilde{k}_m(x,x') \\
&= \sum_{j=1}^{m} S\left( \frac{\pi\,j}{2L} \right) \frac{1}{L}\, \sin\left( \frac{\pi\,j\,(x+L)}{2L} \right) \sin\left( \frac{\pi\,j\,(x'+L)}{2L} \right).
\end{aligned}
$$

(104)

We start by showing the convergence by growing the domain and therefore first consider an approximation with an infinite number of terms $m = \infty$:

$$
\begin{aligned}
&\widetilde{k}_\infty(x, x') \\
&= \sum_{j=1}^\infty S\left(\frac{\pi j}{2L}\right) \frac{1}{L} \sin\left(\frac{\pi j (x + L)}{2L}\right) \sin\left(\frac{\pi j (x' + L)}{2L}\right).
\end{aligned}
$$
(105)

For that purpose we rewrite the summation above in (105) as

$$
\begin{aligned}
&\sum_{j=1}^\infty S\left(\frac{\pi j}{2L}\right) \frac{1}{L} \sin\left(\frac{\pi j (x + L)}{2L}\right) \sin\left(\frac{\pi j (x' + L)}{2L}\right) \\
&= \sum_{j=1}^\infty S\left(\frac{\pi j}{2L}\right) \cos\left(\frac{\pi j (x - x')}{2L}\right) \frac{1}{2L} \\
&\quad - \frac{1}{2L} \sum_{j=1}^\infty \left[ S\left(\frac{\pi 2j}{2L}\right) - S\left(\frac{\pi (2j-1)}{2L}\right) \right] \cos\left(\frac{\pi 2j (x + x')}{2L}\right) \\
&\quad - \frac{1}{2L} \sum_{j=1}^\infty S\left(\frac{\pi (2j-1)}{2L}\right) \left[ \cos\left(\frac{\pi 2j (x + x')}{2L}\right) \right. \\
&\qquad\qquad \left. - \cos\left(\frac{\pi (2j-1)(x + x')}{2L}\right) \right].
\end{aligned}
$$
(106)

and consider the three summations above separately. The analysis of them is done in the next three lemmas.

**Lemma 12.** *Assume that on $\omega \geq 0$ we have $S(\omega) \leq B < \infty$ and $\int_0^\infty S(w)\, d\omega = A < \infty$, and on $\omega > 0$ $S(\omega)$ has a bounded derivative $|S'(\omega)| \leq D < \infty$ and that $\int_0^\infty |S'(\omega)|\, d\omega = C < \infty$. Then there exists a constant $D_2$ such that for all $x, x' \in [-\widetilde{L}, \widetilde{L}]$ we have*

$$
\begin{aligned}
&\left| \sum_{j=1}^\infty S\left(\frac{\pi j}{2L}\right) \cos\left(\frac{\pi j (x - x')}{2L}\right) \frac{1}{2L} \right. \\
&\quad \left. - \frac{1}{\pi} \int_0^\infty S(\omega) \cos(\omega (x - x'))\, d\omega \right| \leq \frac{D_2}{L}.
\end{aligned}
$$
(107)

*Proof.* By using Lemma 9 with $\Delta = \frac{\pi}{2L}$, $f(\omega) = \frac{1}{\pi} S(\omega) \cos(\omega (x - x'))\, d\omega$, $m = 0$, and $\alpha = 0$ as well as the assumptions on $S(\omega)$ and boundedness of sine and cosine we

get that

$$
\begin{aligned}
&\left| \sum_{j=1}^\infty S\left(\frac{\pi j}{2L}\right) \cos\left(\frac{\pi j (x - x')}{2L}\right) \frac{1}{2L} \right. \\
&\quad \left. - \frac{1}{\pi} \int_0^\infty S(\omega) \cos(\omega (x - x'))\, d\omega \right| \\
&\leq \frac{1}{\pi} \int \left| S'(w) \cos(\omega (x - x')) \right. \\
&\qquad \left. - S(w)(x - x') \sin(\omega (x - x')) \right|\, d\omega\, \frac{\pi}{2L} \\
&\leq \frac{1}{2L} \int \left| S'(w) \cos(\omega (x - x')) \right|\, d\omega \\
&\quad + \frac{1}{2L} \int \left| S(w)(x - x') \sin(\omega (x - x')) \right|\, d\omega \\
&\leq \frac{1}{2L} \int \left| S'(w) \right|\, d\omega + \frac{|x - x'|}{2L} \int \left| S(w) \right|\, d\omega \\
&\leq \frac{1}{2L} \int \left| S'(w) \right|\, d\omega + \frac{\widetilde{L}}{L} \int \left| S(w) \right|\, d\omega \\
&\leq \frac{1}{2L} C + \frac{\widetilde{L}}{L} A,
\end{aligned}
$$
(108)

which gives the result with $D_2 = \frac{C}{2} + \widetilde{L} A$. $\qquad\square$

**Lemma 13.** *Assume that for $\omega \geq 0$, $S(\omega)$ is a bounded integrable function with a bounded derivative on $\omega > 0$ such that $\int_0^\infty |S'(\omega)|\, d\omega = C < \infty$, then there exists a constant $D_3$ such that*

$$
\begin{aligned}
&\left| \frac{1}{2L} \sum_{j=1}^\infty \left[ S\left(\frac{\pi 2j}{2L}\right) - S\left(\frac{\pi (2j-1)}{2L}\right) \right] \cos\left(\frac{\pi 2j (x + x')}{2L}\right) \right| \\
&\leq \frac{D_3}{L}.
\end{aligned}
$$
(109)

*Proof.* The result follows by using Lemma 10 with $\Delta = \frac{\pi}{L}$, $\alpha = 1/2$, $\beta = 0$, $f(\omega) = S(\omega)$, and $g(\omega) = \cos(\omega (x + x'))$ and by recalling that $|\cos(\omega (x + x'))| \leq 1$, which gives the constant $D_3 = \frac{C}{2}$. $\qquad\square$

**Lemma 14.** *Assume that for $\omega \geq 0$, $S(\omega)$ is a bounded positive integrable function with bounded derivative on $\omega > 0$ such that $\int_0^\infty S(\omega)\, d\omega = A < \infty$ and $\int_0^\infty |S'(\omega)|\, d\omega = C < \infty$. Then there exists a constant $D_4$ such that*

$$
\begin{aligned}
&\left| \frac{1}{2L} \sum_{j=1}^\infty S\left(\frac{\pi (2j-1)}{2L}\right) \left[ \cos\left(\frac{\pi 2j (x + x')}{2L}\right) \right.\right. \\
&\qquad \left.\left. - \cos\left(\frac{\pi (2j-1)(x + x')}{2L}\right) \right] \right| \leq \frac{D_4}{L}.
\end{aligned}
$$

*Proof.* By using Lemma 11 with $\Delta = \frac{\pi}{L}$, $\alpha = 1/2$, $\beta = 1/2$, $f(\omega) = S(\omega)$, and $g(\omega) = \cos(\omega (x + x'))$ we get

$$
\begin{aligned}
&\left| \frac{1}{2L} \sum_{j=1}^\infty S\left(\frac{\pi (2j-1)}{2L}\right) \left[ \cos\left(\frac{\pi 2j (x + x')}{2L}\right) \right.\right. \\
&\qquad \left.\left. - \cos\left(\frac{\pi (2j-1)(x + x')}{2L}\right) \right] \right| \\
&\leq \frac{(A + C) D'}{2L},
\end{aligned}
$$

where $D'$ is an upper bound for $|(x + x') \sin(\omega (x + x'))|$. We can now select $D' = 2\widetilde{L}$, which gives $D_4 = (A + C) \widetilde{L}$. $\qquad\square$

Next, we combine the above lemmas to get the following result.

**Lemma 15.** *Let the assumptions of Lemmas 12, 13, and 14 be satisfied. Then there exists a constant $D_1$ such that for all $x, x' \in [-\widetilde{L}, \widetilde{L}]$ we have*

$$
\left| \sum_{j=1}^{\infty} S\left(\frac{\pi j}{2L}\right) \frac{1}{L} \sin\left(\frac{\pi j (x+L)}{2L}\right) \sin\left(\frac{\pi j (x'+L)}{2L}\right) \right.
$$
$$
\left. - \frac{1}{\pi} \int_0^{\infty} S(\omega) \cos(\omega (x - x')) \, d\omega \right| \le \frac{D_1}{L}. \quad (110)
$$

*That is,*

$$
\left| \widetilde{k}_\infty(x, x') - k(x, x') \right| \le \frac{D_1}{L}, \quad \text{for } x, x' \in [-\widetilde{L}, \widetilde{L}]. \quad (111)
$$

*Furthermore, the explicit expression for the constant is given as*

$$
D_1 = C + (2A + C)\,\widetilde{L}. \quad (112)
$$

*Proof.* Using triangle inequality to the differece of (106) and $\frac{1}{\pi} \int_0^{\infty} S(\omega) \cos(\omega (x - x')) \, d\omega$ along with Lemmas 12, 13, and 14 gives

$$
\left| \sum_{j=1}^{\infty} S\left(\frac{\pi j}{2L}\right) \cos\left(\frac{\pi j (x - x')}{2L}\right) \frac{1}{2L} \right.
$$
$$
- \frac{1}{2L} \sum_{j=1}^{\infty} \left[ S\left(\frac{\pi 2j}{2L}\right) - S\left(\frac{\pi (2j-1)}{2L}\right) \right] \cos\left(\frac{\pi 2j (x+x')}{2L}\right)
$$
$$
- \frac{1}{2L} \sum_{j=1}^{\infty} S\left(\frac{\pi (2j-1)}{2L}\right) \left[ \cos\left(\frac{\pi 2j (x+x')}{2L}\right) \right.
$$
$$
\left. - \cos\left(\frac{\pi (2j-1)(x+x')}{2L}\right) \right]
$$
$$
\left. - \frac{1}{\pi} \int_0^{\infty} S(\omega) \cos(\omega (x - x')) \, d\omega \right|
$$
$$
\le \left| \sum_{j=1}^{\infty} S\left(\frac{\pi j}{2L}\right) \cos\left(\frac{\pi j (x - x')}{2L}\right) \frac{1}{2L} \right.
$$
$$
\left. - \frac{1}{\pi} \int_0^{\infty} S(\omega) \cos(\omega (x - x')) \, d\omega \right|
$$
$$
+ \left| \frac{1}{2L} \sum_{j=1}^{\infty} \left[ S\left(\frac{\pi 2j}{2L}\right) - S\left(\frac{\pi (2j-1)}{2L}\right) \right] \cos\left(\frac{\pi 2j (x+x')}{2L}\right) \right|
$$
$$
+ \left| \frac{1}{2L} \sum_{j=1}^{\infty} S\left(\frac{\pi (2j-1)}{2L}\right) \left[ \cos\left(\frac{\pi 2j (x+x')}{2L}\right) \right. \right.
$$
$$
\left. \left. - \cos\left(\frac{\pi (2j-1)(x+x')}{2L}\right) \right] \right|
$$
$$
\le \frac{D_2}{L} + \frac{D_3}{L} + \frac{D_4}{L} = \frac{D_1}{L}, \quad (113)
$$

where the explicit values for the costants can be found in the proofs of the lemmas.  □

Let us now consider what happens when we replace the infinite sum approximation with a finite $m$ number of terms. We are now interested in

$$
\widetilde{k}_\infty(x, x') - \widetilde{k}_m(x, x')
$$
$$
= \sum_{j=m+1}^{\infty} S\left(\frac{\pi j}{2L}\right) \frac{1}{L} \sin\left(\frac{\pi j (x+L)}{2L}\right) \sin\left(\frac{\pi j (x'+L)}{2L}\right).
$$

$$
(114)
$$

**Lemma 16.** *Assume that on $\omega \ge 0$, $S(\omega)$ is bounded and integrable, on $\omega > 0$ it has a bounded derivative, and that $\int_0^{\infty} |S'(\omega)| \, d\omega = C < \infty$. Then there exists a constant $D_5$ such that for all $x, x' \in [-\widetilde{L}, \widetilde{L}]$ we have*

$$
\left| \widetilde{k}_\infty(x, x') - \widetilde{k}_m(x, x') \right| \le \frac{D_5}{L} + \frac{2}{\pi} \int_{\frac{\pi m}{2L}}^{\infty} S(\omega) \, d\omega. \quad (115)
$$

*Proof.* Because the sinusoidals are bounded by unity, we get

$$
\left| \sum_{j=m+1}^{\infty} S\left(\frac{\pi j}{2L}\right) \frac{1}{L} \sin\left(\frac{\pi j (x+L)}{2L}\right) \sin\left(\frac{\pi j (x'+L)}{2L}\right) \right|
$$
$$
\le \left| \sum_{j=m+1}^{\infty} S\left(\frac{\pi j}{2L}\right) \frac{1}{L} \right|. \quad (116)
$$

For the right hand side we can now use Lemma 9 with $f(\omega) = \frac{2}{\pi} S(\omega)$ and $\Delta = \frac{\pi}{2L}$, which gives

$$
\left| \sum_{j=m+1}^{\infty} S\left(\frac{\pi j}{2L}\right) \frac{1}{L} - \frac{2}{\pi} \int_{\frac{\pi m}{2L}}^{\infty} S(\omega) \, d\omega \right| \le C \frac{\pi}{2L} = \frac{D_5}{L}. \quad (117)
$$

Hence by the triangle inequality we get

$$
\left| \sum_{j=m+1}^{\infty} S\left(\frac{\pi j}{2L}\right) \frac{1}{L} \right|
$$
$$
= \left| \sum_{j=m+1}^{\infty} S\left(\frac{\pi j}{2L}\right) \frac{1}{L} - \frac{2}{\pi} \int_{\frac{\pi m}{2L}}^{\infty} S(\omega) \, d\omega + \frac{2}{\pi} \int_{\frac{\pi m}{2L}}^{\infty} S(\omega) \, d\omega \right|
$$
$$
\le \left| \sum_{j=m+1}^{\infty} S\left(\frac{\pi j}{2L}\right) \frac{1}{L} - \frac{2}{\pi} \int_{\frac{\pi m}{2L}}^{\infty} S(\omega) \, d\omega \right| + \frac{2}{\pi} \int_{\frac{\pi m}{2L}}^{\infty} S(\omega) \, d\omega
$$
$$
\le \frac{D_5}{L} + \frac{2}{\pi} \int_{\frac{\pi m}{2L}}^{\infty} S(\omega) \, d\omega \quad (118)
$$

and thus the result follows.  □

**Remark 17.** *We can also obtain a bit more defined bound by not using an m-independent bound for forming $D_5$, which under the assumptions of Lemma 16 gives*

$$
\left| \widetilde{k}_\infty(x, x') - \widetilde{k}_m(x, x') \right|
$$
$$
\le \frac{\pi}{2L} \int_{\frac{\pi m}{2L}}^{\infty} |S'(\omega)| \, d\omega + \frac{2}{\pi} \int_{\frac{\pi m}{2L}}^{\infty} S(\omega) \, d\omega. \quad (119)
$$

The lemmas presented in this section can now be combined to a proof of the one-dimensional convergence theorem as follows:

*Proof of Theorem 1.* The first result follows by combining Lemmas 15 and 16 via the triangle inequality. Because our assumptions imply that

$$
\lim_{x \to \infty} \int_x^{\infty} S(\omega) \, d\omega = 0, \quad (120)
$$

for any fixed $L$ we have

$$
\lim_{m \to \infty} \left[ \frac{E}{L} + \frac{2}{\pi} \int_{\frac{\pi m}{2L}}^{\infty} S(\omega) \, d\omega \right] \to \frac{E}{L}. \quad (121)
$$

If we now take the limit $L \to \infty$, the second result in the theorem follows.  □

## A.3 Proof of Theorem 4

When $\mathbf{x} \in \mathbb{R}^d$, the Wiener–Khinchin identity and symmetry of the spectral density imply that

$$
k(\mathbf{x}, \mathbf{x}') = \frac{1}{(2\pi)^d} \int_{\mathbb{R}^d} S(\omega) \exp(-i\,\omega^\mathsf{T}(\mathbf{x} - \mathbf{x}'))\, d\omega
$$
$$
= \frac{1}{\pi^d} \int_0^\infty \cdots \int_0^\infty S(\omega) \prod_{k=1}^d \cos(\omega_k\,(x_k - x_k'))\, d\omega_1 \cdots d\omega_d.
$$
(122)

The $m = \hat{m}^d$ term approximation now has the form

$$
\widetilde{k}_m(\mathbf{x}, \mathbf{x}') = \sum_{j_1,\ldots,j_d=1}^{\hat{m}} S\left(\frac{\pi\,j_1}{2L_1}, \ldots, \frac{\pi\,j_d}{2L_d}\right)
$$
$$
\times \prod_{k=1}^d \frac{1}{L_k} \sin\left(\frac{\pi\,j_k\,(x_k + L_k)}{2L_k}\right) \sin\left(\frac{\pi\,j_k\,(x_k' + L_k)}{2L_k}\right).
$$
(123)

As in the one-dimensional problem we start by considering the case where $\hat{m} = \infty$.

**Lemma 18.** *Let the assumptions of Lemma 15 be satisfied for each $\omega_j \mapsto S(\omega_1, \ldots, \omega_d)$ separately. Then there exists a constant $D_1$ such that for all $\mathbf{x}, \mathbf{x}' \in [-\widetilde{L}, \widetilde{L}]^d$ we have*

$$
\left| \sum_{j_1,\ldots,j_d=1}^\infty S\left(\frac{\pi\,j_1}{2L_1}, \ldots, \frac{\pi\,j_d}{2L_d}\right) \right.
$$
$$
\times \prod_{k=1}^d \frac{1}{L_k} \sin\left(\frac{\pi\,j_k\,(x_k + L_k)}{2L_k}\right) \sin\left(\frac{\pi\,j_k\,(x_k' + L_k)}{2L_k}\right)
$$
$$
\left. - \frac{1}{\pi^d} \int_0^\infty \cdots \int_0^\infty S(\omega) \prod_{k=1}^d \cos(\omega_k\,(x - x'))\, d\omega_1 \cdots d\omega_d \right|
$$
$$
\leq D_1 \sum_{k=1}^d \frac{1}{L_k} \leq \frac{D_1\,d}{L},
$$
(124)

*where $L = \min_k L_k$. That is, for all $\mathbf{x}, \mathbf{x}' \in [-\widetilde{L}, \widetilde{L}]^d$*

$$
\left| \widetilde{k}_\infty(\mathbf{x}, \mathbf{x}') - k(\mathbf{x}, \mathbf{x}') \right| \leq D_1 \sum_{k=1}^d \frac{1}{L_k} \leq \frac{D_1\,d}{L}.
$$
(125)

*Proof.* We can separate the summation over $j_1$ as follows:

$$
\sum_{j_2,\ldots,j_d=1}^\infty \left[ \sum_{j_1=1}^\infty S\left(\frac{\pi\,j_1}{2L_1}, \ldots, \frac{\pi\,j_d}{2L_d}\right) \frac{1}{L_1} \right.
$$
$$
\left. \times \sin\left(\frac{\pi\,j_1\,(x_1 + L_1)}{2L_1}\right) \sin\left(\frac{\pi\,j_1\,(x_1' + L_1)}{2L_1}\right) \right]
$$
$$
\times \prod_{k=2}^d \frac{1}{L_k} \sin\left(\frac{\pi\,j_k\,(x_k + L_k)}{2L_k}\right) \sin\left(\frac{\pi\,j_k\,(x_k' + L_k)}{2L_k}\right).
$$
(126)

By Lemma 15 there now exists a constant $D_{1,1}$ such that

$$
\left| \sum_{j_1=1}^\infty S\left(\frac{\pi\,j_1}{2L_1}, \ldots, \frac{\pi\,j_d}{2L_d}\right) \frac{1}{L_1} \right.
$$
$$
\times \sin\left(\frac{\pi\,j_1\,(x_1 + L_1)}{2L_1}\right) \sin\left(\frac{\pi\,j_1\,(x_1' + L_1)}{2L_1}\right)
$$
$$
\left. - \frac{1}{\pi} \int_0^\infty S\left(\omega_1, \frac{\pi\,j_2}{2L_2}, \ldots, \frac{\pi\,j_d}{2L_d}\right) \cos(\omega_1\,(x_1 - x_1'))\, d\omega_1 \right|
$$
$$
\leq \frac{D_{1,1}}{L_1}.
$$
(127)

The triangle inequality then gives

$$
\left| \sum_{j_1,\ldots,j_d=1}^\infty S\left(\frac{\pi\,j_1}{2L_1}, \ldots, \frac{\pi\,j_d}{2L_d}\right) \right.
$$
$$
\times \prod_{k=1}^d \frac{1}{L_k} \sin\left(\frac{\pi\,j_k\,(x_k + L_k)}{2L_k}\right) \sin\left(\frac{\pi\,j_k\,(x_k' + L_k)}{2L_k}\right)
$$
$$
\left. - \frac{1}{\pi^d} \int_0^\infty \cdots \int_0^\infty S(\omega) \prod_{k=1}^d \cos(\omega_j\,(x_k - x_k'))\, d\omega_1 \cdots d\omega_d \right|
$$
$$
\leq \frac{D_{1,1}}{L_1} + \left| \frac{1}{\pi} \sum_{j_2,\ldots,j_d=1}^\infty \int_0^\infty S\left(\omega_1, \frac{\pi\,j_2}{2L_2}, \ldots, \frac{\pi\,j_d}{2L_d}\right) \right.
$$
$$
\times \cos(\omega_1\,(x_1 - x_1'))\, d\omega_1
$$
$$
\times \prod_{k=2}^d \frac{1}{L_k} \sin\left(\frac{\pi\,j_k\,(x_k + L_k)}{2L_k}\right) \sin\left(\frac{\pi\,j_k\,(x_k' + L_k)}{2L_k}\right)
$$
$$
\left. - \frac{1}{\pi^d} \int_0^\infty \cdots \int_0^\infty S(\omega) \prod_{k=1}^d \cos(\omega_k\,(x_k - x_k'))\, d\omega_1 \cdots d\omega_d \right|.
$$
(128)

We can now similarly bound with respect to the summations over $j_2, \ldots, j_d$ which leads to a bound of the form $\frac{D_{1,1}}{L_1} + \cdots + \frac{D_{1,d}}{L_d}$. Taking $D_1 = \max_k D_{1,k}$ leads to the desired result. $\square$

Now we can consider what happens in the finite truncation of the series. That is, we analyze the following residual sum

$$
\widetilde{k}_\infty(\mathbf{x}, \mathbf{x}') - \widetilde{k}_m(\mathbf{x}, \mathbf{x}')
$$
$$
= \sum_{j_1,\ldots,j_d=\hat{m}+1}^\infty S\left(\frac{\pi\,j_1}{2L_1}, \ldots, \frac{\pi\,j_d}{2L_d}\right)
$$
$$
\times \prod_{k=1}^d \frac{1}{L_k} \sin\left(\frac{\pi\,j_k\,(x_k + L_k)}{2L_k}\right) \sin\left(\frac{\pi\,j_k\,(x_k' + L_k)}{2L_k}\right).
$$
(129)

**Lemma 19.** *Let assumptions of Lemma 16 be satisfied for each $\omega_j \mapsto S(\omega_1, \ldots, \omega_d)$. There exists a constant $D_2$ such that for all $\mathbf{x}, \mathbf{x}' \in [-\widetilde{L}, \widetilde{L}]^d$ we have*

$$
\left| \widetilde{k}_\infty(\mathbf{x}, \mathbf{x}') - \widetilde{k}_m(\mathbf{x}, \mathbf{x}') \right| \leq \frac{D_2\,d}{L} + \frac{1}{\pi^d} \int_{\|\omega\| \geq \frac{\pi\,\hat{m}}{2L}} S(\omega)\, d\omega,
$$
(130)

*where $L = \min_k L_k$.*

*Proof.* We can write the following bound

$$
\left| \sum_{j_1,\ldots,j_d=\hat{m}+1}^{\infty} S\left(\frac{\pi j_1}{2L_1},\ldots,\frac{\pi j_d}{2L_d}\right) \right.
$$

$$
\left. \times \prod_{k=1}^{d} \frac{1}{L_k} \sin\left(\frac{\pi j_k (x_k+L_k)}{2L_k}\right) \sin\left(\frac{\pi j_k (x'_k+L_k)}{2L_k}\right) \right|
$$

$$
\leq \left| \sum_{j_1,\ldots,j_d=\hat{m}+1}^{\infty} S\left(\frac{\pi j_1}{2L_1},\ldots,\frac{\pi j_d}{2L_d}\right) \prod_{k=1}^{d}\frac{1}{L_k} \right|. \tag{131}
$$

We can now use Lemma 9 with $f(\omega_1) = \frac{2}{\pi} S\left(\omega_1,\frac{\pi j_2}{2L_2},\ldots,\frac{\pi j_d}{2L_d}\right)$ and $\Delta = \frac{\pi}{2L_1}$, which gives

$$
\left| \sum_{j_1,\ldots,j_d=\hat{m}+1}^{\infty} S\left(\frac{\pi j_1}{2L_1},\ldots,\frac{\pi j_d}{2L_d}\right) \prod_{k=1}^{d}\frac{1}{L_k} \right.
$$

$$
\left. -\frac{2}{\pi} \sum_{j_2,\ldots,j_d=\hat{m}+1}^{\infty} \int_{\frac{\pi\hat{m}}{2L_1}}^{\infty} S\left(\omega_1,\frac{\pi j_2}{2L_2},\ldots,\frac{\pi j_d}{2L_d}\right) \mathrm{d}\omega_1 \prod_{k=2}^{d}\frac{1}{L_k} \right|
$$

$$
\leq \frac{D_{2,1}}{L_1}. \tag{132}
$$

Using a similar argument again, we get

$$
\left| \frac{2}{\pi} \sum_{j_2,\ldots,j_d=\hat{m}+1}^{\infty} \int_{\frac{\pi\hat{m}}{2L_1}}^{\infty} S\left(\omega_1,\frac{\pi j_2}{2L_2},\ldots,\frac{\pi j_d}{2L_d}\right) \mathrm{d}\omega_1 \prod_{k=2}^{d}\frac{1}{L_k} \right.
$$

$$
-\frac{2^2}{\pi^2} \sum_{j_3,\ldots,j_d=\hat{m}+1}^{\infty} \int_{\frac{\pi\hat{m}}{2L_1}}^{\infty}\int_{\frac{\pi\hat{m}}{2L_2}}^{\infty} S\left(\omega_1,\omega_2,\frac{\pi j_3}{2L_3},\ldots,\frac{\pi j_d}{2L_d}\right) \mathrm{d}\omega_1\, \mathrm{d}\omega_2
$$

$$
\left. \prod_{k=3}^{d}\frac{1}{L_k} \right| \leq \frac{D_{2,2}}{L_2}. \tag{133}
$$

After repeating this for all the indexes, by forming a telescoping sum of the terms and applying the triangle inequality then gives

$$
\left| \sum_{j_1,\ldots,j_d=\hat{m}+1}^{\infty} S\left(\frac{\pi j_1}{2L_1},\ldots,\frac{\pi j_d}{2L_d}\right) \prod_{k=1}^{d}\frac{1}{L_k} \right.
$$

$$
\left. -\left(\frac{2}{\pi}\right)^d \int_{\frac{\pi\hat{m}}{2L_1}}^{\infty}\cdots\int_{\frac{\pi\hat{m}}{2L_d}}^{\infty} S(\omega_1,\ldots,\omega_d)\,\mathrm{d}\omega_1\cdots\mathrm{d}\omega_d \right| \leq \sum_{k=1}^{d}\frac{D_{2,k}}{L_k}. \tag{134}
$$

Applying the triangle inequality again gives

$$
\left| \sum_{j_1,\ldots,j_d=\hat{m}+1}^{\infty} S\left(\frac{\pi j_1}{2L_1},\ldots,\frac{\pi j_d}{2L_d}\right) \prod_{k=1}^{d}\frac{1}{L_k} \right|
$$

$$
\leq \sum_{k=1}^{d}\frac{D_{2,k}}{L_k} + \left(\frac{2}{\pi}\right)^d \int_{\frac{\pi\hat{m}}{2L_1}}^{\infty}\cdots\int_{\frac{\pi\hat{m}}{2L_d}}^{\infty} S(\omega_1,\ldots,\omega_d)\,\mathrm{d}\omega_1\cdots\mathrm{d}\omega_d. \tag{135}
$$

By interpreting the latter integral as being over the positive exterior of a rectangular hypercuboid and bounding it by a integral over exterior of a hypersphere which fits inside the cuboid, we can bound the expression by

$$
\sum_{k=1}^{d}\frac{D_{2,k}}{L_k} + \frac{1}{\pi^d}\int_{\|\omega\|\geq\frac{\pi\hat{m}}{2L}} S(\omega)\,\mathrm{d}\omega. \tag{136}
$$

The first term can be further bounded by replacing $L_k$s with their minimum $L$ and by defining $D_2 = \max D_{2,k}$ which is $d$ times the maximum of $D_{2,k}$. This leads to the final form of the result.   □

**Remark 20.** *Note that analogously to Remark 17 we could tighten the bound for $D_2$ by letting it depend on $\hat{m}$.*

*Proof of Theorem 4.* Analogous to the one-dimensional case. That is, we combine the results of the above lemmas using the triangle inequality.   □

## References

Adam V, Hensman J, Sahani M (2016) Scalable transformed additive signal decomposition by non-conjugate Gaussian process inference. In: IEEE International Workshop on Machine Learning for Signal Processing (MLSP)

Adler RJ (1981) The geometry of random fields, vol 62. Siam

Akhiezer NI, Glazman IM (1993) Theory of Linear Operators in Hilbert Space. Dover, New York

Álvarez MA, Luengo D, Lawrence ND (2013) Linear latent force models using Gaussian processes. IEEE Transactions on Pattern Analysis and Machine Intelligence 35(11):2693–2705

Bach F (2013) Sharp analysis of low-rank kernel matrix approximations. In: Proceedings of the 26th Annual Conference on Learning Theory (COLT), PMLR, Princeton, NJ, USA, Proceedings of Machine Learning Research, vol 30, pp 185–209

Baker CTH (1977) The Numerical Treatment of Integral Equations. Clarendon press, Oxford

Boutsidis C, Gittens A (2013) Improved matrix algorithms via the subsampled randomized Hadamard transform. SIAM Journal on Matrix Analysis and Applications 34(3):1301–1340

Brooks S, Gelman A, Jones GL, Meng XL (2011) Handbook of Markov Chain Monte Carlo. Chapman & Hall/CRC

Bui TD, Yan J, Turner RE (2017) A unifying framework for Gaussian process pseudo-point approximations using power expectation propagation. Journal of Machine Learning Research 18(104):1–72

Chalupka K, Williams CKI, Murray I (2013) A framework for evaluating approximation methods for Gaussian process regression. Journal of Machine Learning Research 14:333–350

Courant R, Hilbert D (2008) Methods of Mathematical Physics, vol 1. Wiley-VCH

Cramér H, Leadbetter MR (2013) Stationary and Related Stochastic Processes: Sample Function Properties and Their Applications. Dover, Mineola, NY

Csató L, Opper M (2002) Sparse online Gaussian processes. Neural Computation 14(3):641–668

Da Prato G, Zabczyk J (1992) Stochastic Equations in Infinite Dimensions, Encyclopedia of Mathematics and its Applications, vol 45. Cambridge University Press

Deisenroth MP, Ng JW (2015) Distributed Gaussian processes. In: International Conference on Machine Learning (ICML), pp 1481–1490

Duane S, Kennedy AD, Pendleton BJ, Roweth D (1987) Hybrid Monte Carlo. Physics Letters B 195(2):216–222

Feller W (1968) An introduction to probability theory and its applications, vol I, 3rd edn. Wiley

Fritz J, Neuweiler I, Nowak W (2009) Application of FFT-based algorithms for large-scale universal kriging problems. Mathematical Geosciences 41(5):509–533

Gal Y, Turner R (2015) Improving the Gaussian process sparse spectrum approximation by representing uncertainty in frequency inputs. In: Proceedings of the 32nd

International Conference on Machine Learning, Proceedings of Machine Learning Research, vol 37, pp 655–664

Gardner J, Pleiss G, Wu R, Weinberger K, Wilson A (2018) Product kernel interpolation for scalable gaussian processes. In: Proceedings of the Twenty-First International Conference on Artificial Intelligence and Statistics (AISTATS), PMLR, Playa Blanca, Lanzarote, Canary Islands, Proceedings of Machine Learning Research, vol 84, pp 1407–1416

Golub GH, Van Loan CF (1996) Matrix Computations, 3rd edn. Johns Hopkins University Press, Baltimore

Harbrecht H, Peters M, Schneider R (2012) On the low-rank approximation by the pivoted Cholesky decomposition. Applied Numerical Mathematics 4(62):428–440

Hensman J, Fusi N, Lawrence ND (2013) Gaussian processes for big data. In: Proceedings of the 29th Conference on Uncertainty in Artificial Intelligence (UAI 2013), pp 282–290

Hensman J, Durrande N, Solin A (2018) Variational Fourier features for Gaussian processes. Journal of Machine Learning Research 8(151):1–52

Izmailov P, Novikov A, Kropotov D (2018) Scalable Gaussian processes with billions of inducing inputs via tensor train decomposition. In: Proceedings of the Twenty-First International Conference on Artificial Intelligence and Statistics (ICML), PMLR, Playa Blanca, Lanzarote, Canary Islands, Proceedings of Machine Learning Research, vol 84, pp 726–735

Kaipio J, Somersalo E (2005) Statistical and Computational Inverse Problems. Springer

Kimeldorf GS, Wahba G (1970) A correspondence between Bayesian estimation on stochastic processes and smoothing by splines. The Annals of Mathematical Statistics pp 495–502

Lázaro-Gredilla M (2010) Sparse Gaussian processes for large-scale machine learning. PhD thesis, Universidad Carlos III de Madrid, Madrid, Spain

Lázaro-Gredilla M, Quiñonero-Candela J, Rasmussen CE, Figueiras-Vidal AR (2010) Sparse spectrum Gaussian process regression. Journal of Machine Learning Research 11:1865–1881

Le Q, Sarlos T, Smola A (2013) Fastfood – Computing Hilbert space expansions in loglinear time. In: Proceedings of the 30th International Conference on Machine Learning (ICML), PMLR, Atlanta, Georgia, USA, Proceedings of Machine Learning Research, vol 28, pp 244–252

Lenk PJ (1991) Towards a practicable Bayesian nonparametric density estimator. Biometrika 78(3):531–543

Lindgren F, Rue H, Lindström J (2011) An explicit link between Gaussian fields and Gaussian Markov random fields: The stochastic partial differential equation approach. Journal of the Royal Statistical Society: Series B (Statistical Methodology) 73(4):423–498

Liu JS (2001) Monte Carlo Strategies in Scientific Computing. Springer, New York

Loève M (1963) Probability Theory, 3rd edn. The University Series in Higher Mathematics, Van Nostrand, Princeton, NJ

Neal RM (2011) MCMC using Hamiltonian dynamics. In: Brooks S, Gelman A, Jones GL, Meng XL (eds) Handbook of Markov Chain Monte Carlo, Chapman & Hall/CRC, chap 5

Opper M, Vivarelli F (1999) General bounds on Bayes errors for regression with Gaussian processes. In: Advances in Neural Information Processing Systems, vol 11, pp 302–308

Paciorek CJ (2007) Bayesian smoothing with Gaussian processes using Fourier basis functions in the spectralGP package. Journal of Statistical Software 19(2):1–38

Quiñonero-Candela J, Rasmussen CE (2005a) Analysis of some methods for reduced rank Gaussian process regression. In: Switching and Learning in Feedback Systems, Lecture Notes in Computer Science, vol 3355, Springer, pp 98–127

Quiñonero-Candela J, Rasmussen CE (2005b) A unifying view of sparse approximate Gaussian process regression. Journal of Machine Learning Research 6:1939–1959

Rasmussen CE, Williams CKI (2006) Gaussian Processes for Machine Learning. The MIT Press

Saatçi Y (2012) Scalable inference for structured Gaussian process models. PhD thesis, University of Cambridge, UK

Samo YLK, Roberts SJ (2016) String and membrane Gaussian processes. Journal of Machine Learning Research 17:1–87

Särkkä S (2011) Linear operators and stochastic partial differential equations in Gaussian process regression. In: Proceedings of ICANN

Särkkä S, Hartikainen J (2012) Infinite-dimensional Kalman filtering approach to spatio-temporal Gaussian process regression. In: Proceedings of the Fifteenth International Conference on Artificial Intelligence and Statistics (AISTATS), JMLR Workshop and Conference Proceedings, vol 22, pp 993–1001

Särkkä S, Piché R (2014) On convergence and accuracy of state-space approximations of squared exponential covariance functions. In: Proceedings of MLSP, pp 1–6

Särkkä S, Solin A, Hartikainen J (2013) Spatiotemporal learning via infinite-dimensional Bayesian filtering and smoothing. IEEE Signal Processing Magazine 30(4):51–61

Seeger M, Williams CKI, Lawrence ND (2003) Fast forward selection to speed up sparse Gaussian process regression. In: Proceedings of the 9th International Workshop on Artificial Intelligence and Statistics (AISTATS)

Showalter RE (2010) Hilbert Space Methods in Partial Differential Equations. Dover Publications

Shubin MA (1987) Pseudodifferential Operators and Spectral Theory. Springer Series in Soviet Mathematics, Springer-Verlag

Smola AJ, Bartlett P (2001) Sparse greedy Gaussian process regression. In: Advances in Neural Information Processing Systems, vol 13

Snelson E, Ghahramani Z (2006) Sparse Gaussian processes using pseudo-inputs. In: Advances in Neural Information Processing Systems, vol 18, pp 1259–1266

Sollich P, Halees A (2002) Learning curves for Gaussian process regression: Approximations and bounds. Neural Computation 14(6):1393–1428

Tarantola A (2004) Inverse Problem Theory and Methods for Model Parameter Estimation. SIAM

Titsias MK (2009) Variational learning of inducing variables in sparse Gaussian processes. In: Proceedings of the 12th International Conference on Artificial Intelligence and Statistics (AISTATS), JMLR Workshop and Conference Proceedings, vol 5, pp 567–574

Tresp V (2000) A Bayesian committee machine. Neural Computation 12(11):2719–2741

Van Trees HL (1968) Detection, Estimation, and Modulation Theory Part I. John Wiley & Sons, New York

Vanhatalo J, Pietiläinen V, Vehtari A (2010) Approximate inference for disease mapping with sparse Gaussian processes. Statistics in Medicine 29(15):1580–1607

Vanhatalo J, Riihimäki J, Hartikainen J, Jylänki P, Tolvanen V, Vehtari A (2013) GPstuff: Bayesian modeling with Gaussian processes. Journal of Machine Learning Research 14:1175–1179

Wahba G (1978) Improper priors, spline smoothing and the problem of guarding against model errors in regression. Journal of the Royal Statistical Society Series B (Methodological) pp 364–372

Wahba G (1990) Spline Models for Observational Data. SIAM

Williams CKI, Seeger M (2000) The effect of the input density distribution on kernel-based classifiers. In: Proceedings of the 17th International Conference on Machine Learning

Williams CKI, Seeger M (2001) Using the Nyström method to speed up kernel machines. In: Advances in Neural Information Processing Systems, vol 13

Wilson A, Nickisch H (2015) Kernel interpolation for scalable structured Gaussian processes (KISS-GP). In: Proceedings of the 32nd International Conference on Machine Learning (ICML), PMLR, Lille, France, Proceedings of Machine Learning Research, vol 37, pp 1775–1784