
Sparse Spectral Bayesian Permanental Process with Generalized Kernel

Jeremy Sellier

Department of Statistical Science
University College London, UK

Petros Dellaportas

Department of Statistical Science
University College London, UK
Department of Statistics,
Univ. of Econ. and Business, Athens, Greece
and The Alan Turing Institute, UK

Abstract

We introduce a novel scheme for Bayesian inference on permanental processes which models the Poisson intensity as the square of a Gaussian process. Combining *generalized kernels* and a Fourier features-based representation of the Gaussian process with a Laplace approximation to the posterior, we achieve a fast and efficient inference that does not require numerical integration over the input space, allows kernel design and scales linearly with the number of events. Our method builds and improves upon the state-of-the-art Laplace Bayesian point process benchmark of Walder and Bishop (2017), demonstrated on both synthetic, real-world temporal and large spatial data sets.

1 INTRODUCTION

The Poisson process is commonly used in a variety of real-world point pattern applications including seismic activity (Gardner and Knopoff, 1974), epidemiology (Diggle, 2003; Banerjee et al., 2003), neuroscience (Cunningham, Yu, Shenoy and Sahani, 2008) and crime incident locations (Grubestic and Mack, 2008; Flaxman et al., 2019). A flexible generalization, incorporating a stochastic Poisson intensity (Cox, 1955) is achieved by placing a non-parametric Gaussian process (GP) prior over the intensity function, resulting in the popular Gaussian Cox process model.

Inference with the Gaussian Cox process model is challenging due to the *doubly-intractable* likelihood, requiring integration of an infinite-dimensional random function over the input domain. Furthermore, a positive transformation must

be applied to the GP prior to ensure the intensity function remains non-negative. Different choices of transformation function are found in the literature. A first classical approach relies on the exponential transformation, resulting in the *log Gaussian Cox process* proposed by Møller et al. (1998). This model usually requires numerical approximation of the integral by discretization over the input space (Møller et al., 1998; Diggle et al., 2013), a computationally-intensive procedure which scales poorly with the dimensionality of the input domain. A second approach uses a sigmoid transformation and an input space augmentation via *thinning*, to construct an exact Markov Chain Monte Carlo sampler (Adams et al., 2009; Gunter et al., 2014), eliminating the need for likelihood integration. In practice however, this approach is computationally intractable for large problems. Other works include the use of the relu (Ko and Seeger, 2016) and softplus (Seeger and Bouchard, 2012; Park et al., 2014) as transformation functions. Finally, Lopez-lopera et al. (2019) introduce a finite approximation where positivity conditions is imposed directly on the GP.

Another approach exploits the so-called *permanental process*, defining the Poisson process intensity in terms of the square of a GP (McCullagh and Møller, 2006; Lloyd et al., 2015). It enables analytical computation of the intensity integral when coupled with a variational inference scheme with inducing points (similar to Titsias 2009) and has received considerable recent attention (Lian et al., 2015; Flaxman et al., 2017; John and Hensman, 2018).

Walder and Bishop (2017) propose a Laplace Bayesian point process (LBPP) method, a fast alternative for the permanental process that relies on the Mercer decomposition of the Gaussian process kernel and a Laplace approximation to the intensity posterior. They show significant speed improvement compared to variational Bayesian inference. Inference based on the Laplace approximation has already been proposed in the context of a Gaussian Cox process by Cunningham, Shenoy and Sahani (2008), Illian et al. (2012), and Flaxman et al. (2015).

However, the tractability properties of the permanental pro-

cess used by Lloyd et al. (2015) and Walder and Bishop (2017) only holds for certain standard types of kernels such as the squared exponential kernel, which encodes restrictive assumptions about the form of the function we are modelling. In general, the choice of kernel determines almost all the generalization properties of a Gaussian process model and profoundly affects its performance on a given task (Rasmussen and Williams, 2005). Approaches have been proposed in recent years to achieve more expressible kernels either by a composition of simple analytical forms (Duvenaud et al., 2011, 2013) or more flexibly through a spectral representation (Lázaro-Gredilla et al., 2010; Wilson and Adams, 2013; Samo and Roberts, 2015a).

In this paper, we build on the LBPP approach of Walder and Bishop (2017), introducing an alternative fast Laplace-based inference exploiting spectral representation of kernels and random Fourier features (RFFs). Our approach, the Sparse Spectral Permanental Process (SSPP), retains the tractability properties of the permanental process, whilst being able to adapt to a broader range of stationary kernels. Furthermore, our method works with *generalized stationary spectral kernels* (Samo and Roberts, 2015a), to our knowledge, the most general class of expressible spectral kernels, that can approximate any stationary kernels to arbitrary precision.

Following John and Hensman (2018), we also include a mean constant to mitigate the effect of *nodal lines* observed for the permanental process, resulting from the non-injective nature of the squared transformation. Our approach shows systematic improvement in accuracy in synthetic and real-world data sets.

2 PRELIMINARIES

2.1 Gaussian Cox Process

A Poisson process (see Daley and Vere-Jones, 2003) models a random sequence of points occurring on a continuous domain X . Inference involves estimation of an intensity function $\lambda(\mathbf{x}) : X \rightarrow \mathbb{R}^+$, that can be interpreted heuristically as the instantaneous probability of occurrence of a point around a location $\mathbf{x} \in X$, i.e. $\lambda(\mathbf{x}) := \lim_{(d\mathbf{x}) \rightarrow 0^+} \mathbb{E}[N(d\mathbf{x})] / (d\mathbf{x})$, where $d\mathbf{x}$ is a small neighborhood around \mathbf{x} with measure $(d\mathbf{x})$ and $N(A)$ is the random number of points within a sub-region $A \subset X$.

We focus our attention on Gaussian Cox process models, where λ is defined as $\lambda(\mathbf{x}) := I(f(\mathbf{x}))$, for a non-negative transformation $I : \mathbb{R} \rightarrow \mathbb{R}^+$ and a function $f \sim GP(0; k(\mathbf{x}; \mathbf{x}^\theta))$, with $k : X \times X \rightarrow \mathbb{R}$ being the positive-definite covariance function for f . Assuming $\mathbf{X} = \{\mathbf{x}_i\}_{i=1}^N$ to be a realization of N observations in X , the sample likelihood is

$$p(\mathbf{X}|\mathbf{f}) = \exp \left(\int_X I(f(\mathbf{x})) d\mathbf{x} - \sum_{i=1}^N I(f(\mathbf{x}_i)) \right) \quad (1)$$

and a latent posterior $p(\mathbf{f}|\mathbf{X})$ is

$$\frac{\int_X \exp \left(\int_X I(f(\mathbf{x})) d\mathbf{x} - \sum_{i=1}^N I(f(\mathbf{x}_i)) \right) p(\mathbf{f}) d\mathbf{f}}{\int_X \exp \left(\int_X I(f(\mathbf{x})) d\mathbf{x} - \sum_{i=1}^N I(f(\mathbf{x}_i)) \right) p(\mathbf{f}) d\mathbf{f}} \quad (2)$$

where \mathbf{f} represents the infinite-dimensional object corresponding to $f(\mathbf{x})$. Equation (2) is often described as “doubly-intractable”. Inference for \mathbf{f} requires evaluating integral terms that cannot be calculated explicitly.

2.2 Permanental Process

The permanental process (McCullagh and Møller, 2006) is obtained by defining the intensity in Equation (1) as the square of a Gaussian process i.e. setting $\lambda(\mathbf{x}) = f(\mathbf{x})^2$. This transformation is advantageous, in that the density of the resulting permanental process is available analytically. To express the density we make use of the spectral decomposition of the Gaussian process covariance function, which we shortly review in this section.

2.2.1 Integral Expression via Mercer Theorem

The Gaussian process covariance function k has a Mercer decomposition on $(X; \cdot)$, if it can be written as

$$k(\mathbf{x}; \mathbf{x}^\theta) = \sum_{i=1}^{\infty} g_i(\mathbf{x}) g_i(\mathbf{x}^\theta) \quad \text{for } \mathbf{x}, \mathbf{x}^\theta \in X \quad (3)$$

where $\{g_i\}_{i=1}^{\infty}$ is a sequence of summable, non-negative, non-increasing *eigenvalues*, and $\{g_i(\cdot)\}_{i=1}^{\infty}$ is a set of mutually-orthogonal, unit-norm *eigenfunctions* with respect to the inner product $\langle u, v \rangle = \int_X u(\mathbf{x})v(\mathbf{x})d(\mathbf{x})$. Then in a similar manner to McCullagh and Møller (2006), Flaxman et al. (2017) and Walder and Bishop (2017), $f(\mathbf{x}) \sim GP(0; k(\mathbf{x}; \mathbf{x}^\theta))$ can be reformulated as an equivalent linear form

$$f(\mathbf{x}) = \sum_{i=0}^{\infty} w_i g_i(\mathbf{x}) \quad (4)$$

where $\mathbf{w} = (w_1; w_2; \dots)^T \sim N(\mathbf{0}; \Sigma)$ and Σ is a diagonal covariance matrix with entries $\Sigma_{ii} = g_i; i = 1; 2; \dots$. Further, it can be shown that $\text{Cov}(f(\mathbf{x}); f(\mathbf{x}^\theta)) = \mathbf{w}^T \mathbf{w} = k(\mathbf{x}; \mathbf{x}^\theta)$, where \mathbf{w} is a vector with entries $w_i; i = 1; 2; \dots$. The integral of the intensity can then be expressed as

$$\begin{aligned} \int_X f(\mathbf{x})^2 d(\mathbf{x}) &= \sum_{i=0}^{\infty} \sum_{j=0}^{\infty} w_i w_j \int_X g_i(\mathbf{x}) g_j(\mathbf{x}) d(\mathbf{x}) \\ &= \sum_{i=0}^{\infty} w_i^2 \end{aligned} \quad (5)$$

2.2.2 Approximate Bayesian Inference

In our case, to make the reformulation of the integral $\int_{\mathcal{S}} f(x)^2 dx$ possible as in Equation (5), the kernel requires an explicit Mercer representation with respect to the Lebesgue measure; this is not available for most choices of kernel. In such cases, the Nyström method can be used to approximate the eigenfunctions and eigenvalues of Equation (3). Both Flaxman et al. (2017) and Walder and Bishop (2017) adopt the Nyström approach in the context of the Permanent Cox process with Gaussian kernel. Walder and Bishop (2017) further propose a Bayesian inference scheme based on a Laplace approximation for a non-GP likelihood. A more detailed review of the Laplace approach together with a new result for the integral in Equation (5) with the Nyström method is presented in part A of the Appendix.

3 MODEL

Motivated in part by the shortcomings of the Nyström approach proposed by Walder and Bishop (2017), we now present an alternative LBBP approach to inference for the permanent process. In contrast to the Mercer approach, it is based on a sparse spectral representation of a GP, exploiting random Fourier features (RFFs, Rahimi and Recht, 2007) for reduced-rank kernel expression. As a result, it provides a tractable expression for the integral of the intensity over the input domain.

Our spectral approach works for any bounded, continuous and shift-invariant kernel $k(x; x^0) := k(x - x^0)$ that satisfies the condition of Bochner's theorem (see Theorem 3.1) and admits a finite dimensional feature space representation or approximation. In contrast, the variational inference approach of Lloyd et al. (2015) and the LBPP with Nyström, yield an analytical integral expression for a limited choice of kernels, like the Gaussian kernel. Furthermore, we are able to adapt our method to generalized stationary spectral kernels (Samo and Roberts, 2015) which generalize two other classes of expressible spectral kernels, sparse spectrum kernels (Lázaro-Gredilla et al., 2010) and mixture spectral kernels (Wilson and Adams, 2013). These two kernels have been proven to be able to approximate any bounded continuous stationary kernels to arbitrary precision.

We also address the issue of nodal lines discussed in John and Hensman (2018). This problem arises since the inverse link function $\psi(\cdot) := f(\cdot)^2$ is not injective, with $f(\cdot)$ producing the same intensity. Therefore, regions of negative and positive f must exhibit zero-crossings, where the intensity is artificially forced to zero, despite the underlying intensity being positive. Following John and Hensman (2018), we add an offset parameter to the intensity function $\psi(\cdot) := (f(\cdot) + \epsilon)^2$ corresponding to an initial value for the prior mean of the GP, to alleviate the problem.

3.1 Sparse Spectral Kernels

In this section, we briefly present two families of spectral kernels, sparse spectrum kernels (Lázaro-Gredilla et al., 2010) and mixture spectral kernels (Wilson and Adams, 2013) that have been proposed in recent years for kernel design. They are both known to be dense in the family of stationary kernels, implying that they can approximate any stationary kernel to an arbitrary precision given sufficient spectral components.

Spectral kernels are constructed via the Bochner's theorem (Bochner, 1932), which states that any bounded, continuous and shift-invariant kernel $k(x; x^0) := k(\cdot)$ with $\cdot = x - x^0$, is the inverse Fourier transform of a bounded positive measure.

Theorem 3.1. (Bochner) An integrable function $k: \mathbb{R}^d \rightarrow \mathbb{R}$ is the covariance function of a weakly stationary mean square continuous random process $\{X_t\}_{t \in \mathbb{R}^d}$ if and only if it can be represented as

$$k(\cdot) = \int_{\mathbb{R}^d} \exp(i\langle \cdot, z \rangle) d\mu(z) \quad (6)$$

where μ is a positive definite measure.

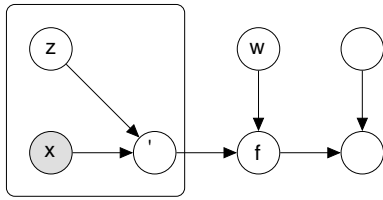
Sparse spectrum kernels can be obtained by setting in Equation (6) to be a positive discrete symmetric measure $\mu = \sum_{k=1}^K \frac{a_k}{2} (\delta_{\omega_k} + \delta_{-\omega_k})$ where $a_k > 0$ and δ_{ω_k} denotes the Dirac measure centred at the point spectral frequency $\omega_k \in \mathbb{R}^d$ for $k = 1, \dots, K$. Note that as such, μ is singular with respect to the Lebesgue measure and does not admit a density. Through Equation (6), we obtained the sparse spectrum kernel spanned by the trigonometric functions $\cos(\langle \cdot, \omega_k \rangle) g_{k=1}^K$. A major challenge is that a direct optimization of the linear coefficients $a_k g_{k=1}^K$ and the frequencies $\omega_k g_{k=1}^K$ often leads to overfitting as illustrated by Lázaro-Gredilla et al. (2010) in the context of GP regression.

Wilson and Adams (2013) consider the case where μ is absolutely continuous with respect to the Lebesgue measure and admits a spectral density $\rho(\cdot)$. In that case $\mu(\cdot)$ and the kernel function k are Fourier duals of each other. Mixture spectral kernels model the spectral density as a mixture of independent Gaussian densities with non-zero mean. Since mixtures of Gaussians are dense in the set of all distribution functions (Plataniotis and Hatzinakos, 2001), the resulting dual of this set is dense in the family of continuous stationary kernels.

Random Fourier Features Random Fourier features (Rahimi and Recht, 2007) are closely related to sparse spectrum kernel. From Bochner's theorem, a kernel function

¹Note that positive finite discrete measures are weakly dense in the space of all positive finite measure (Hu and Papageorgiou, 2013).

Figure 1: Illustration of the model, where the arrows directions suggest directions of influence.



can be rewritten as

$$k(x, x^0) = \frac{1}{r} \prod_{k=1}^r \exp(i z_k^>(x - x^0)) \quad (7)$$

where z_1, \dots, z_r in \mathbb{R}^d are independent samples from the distribution with density $S(\cdot)$, for some integer $r > 0$ and $r > 0$. Here, we have assumed that Equation (6) is absolutely continuous with respect to the Lebesgue measure and has a spectral density $S(\cdot)$. Equation (6) is then approximated using Monte Carlo integration. We also treat the scale parameter of the kernel function separately for convenience.

We thus obtain a kernel approximation

$$k(x, x^0) \approx \prod_{k=1}^r g^{(k)}(x) g^{(k)}(x^0) \quad (8)$$

where $g^{(r)}$ is an explicit feature mapping $g^{(r)} : X \rightarrow \mathbb{R}^r$ such that

$$g^{(r)}(x) = \left[\exp(i z_1^>x); \dots; \exp(i z_r^>x) \right]^> \quad (9)$$

We may obtain a $2r$ -sized real-valued mapping that satisfies Equation (7) using

$$g^{(r)}(x) = \left[\cos(z_1^>x); \dots; \cos(z_r^>x); \sin(z_1^>x); \dots; \sin(z_r^>x) \right]^> \quad (10)$$

where $z \sim S(z)$.

The derivation of Equation (10) is provided in part C of the Appendix. Some common kernels and their corresponding spectral densities $S(\cdot)$ are presented in part B of the Appendix. Inspecting Equation (10), we see RFF as a special case of sparse spectrum kernels where the frequencies z_k are sampled at random from some distribution rather than optimised. However, RFF methods do not address the need for explicitly learning the spectral measure from the data.

3.2 Generalized Stationary Kernels

One advantage of our method is that it can work with the generalized stationary kernels (Samo and Roberts, 2015) that are dense in the family of stationary kernels and admits sparse spectrum kernels and mixture spectral kernels as special cases. Generalized kernels can also account for different

degree of differentiability of the latent function. Sparse spectrum kernels and mixture spectral kernels are more limited in a sense that, when used as covariance functions, they yield infinite differentiability of the corresponding stochastic process, which might be unrealistic for certain learning tasks (Stein, 1999).

Definition 3.2. (Generalized stationary kernel) Let g be a stationary kernel $g : \mathbb{R}^d \rightarrow \mathbb{R}$ such that $g(0) = 1$. A generalized kernel k_{GS} with $K \geq N^+$ components takes the form

$$k_{GS}(x, x^0) = \prod_{k=1}^K g_k(x - x^0) \cos(\lambda_k^>(x - x^0)) \quad (11)$$

where $\lambda_k \in \mathbb{R}^d$, $\lambda_k \in \mathbb{R}^{+d}$, $\lambda_k > 0$ for $k = 1, \dots, K$ and \odot denotes the element-wise Hadamard product.

The parameters λ_k are used as inverse input scales. When λ_k are set to zero, we retrieve the sparse spectrum kernels. The spectral mixture kernel corresponds to a special case where $g(x) = \exp(-\|x\|_2^2) = \frac{1}{2}$.

The degree of smoothness of a zero-mean GP with kernel k_{GS} is determined by the kernel. Samo (2017) proposes learning the differentiability of the underlying latent function, by setting g to be a Matérn kernel with different parameter values from $\frac{1}{2} + i$; $i = 0, 1, 2$. The case of $i = 0$ corresponds to continuity and the case of $i = 1$ to 1 times differentiability.

Finite-dimensional Feature Space Approximation For our methodology, we are interested in having a reduced rank representation for k_{GS} similar to Equation (10). Any consistent RFF approximation of g in Equation (11) such that $g(x, x^0) \approx \prod_{k=1}^K g_k(x) g_k(x^0)$ where g_k is an explicit feature mapping $g_k : X \rightarrow \mathbb{R}^r$, would result in a finite-dimensional feature space approximation k_{GS} (Samo, 2017). In that case,

$$k_{GS}(x, x^0) \approx \prod_{k=1}^K h_k(x) h_k(x^0)$$

$$h_k(x) = \left[g_k^{(r)}(x) \cos(\lambda_k^>x); g_k^{(r)}(x) \sin(\lambda_k^>x) \right]^> \quad (12)$$

for $k = 1, \dots, K$, where \odot denotes the Kronecker product.

To be consistent with previous notations, we define $k_{GS}(x, x^0) \approx \prod_{k=1}^K h_k(x) h_k(x^0)$ where $h_k^{(r)}(x)$ is a feature mapping $h_k^{(r)}(x) : X \rightarrow \mathbb{R}^{2rK}$ satisfying

$$h_k^{(r)}(x) = \left[h_{k,1}^{(r)}(x); \dots; h_{k,2r}^{(r)}(x) \right]^> \quad (13)$$

3.3 Sparse Spectral Permanental Process (SSPP)

Using RFFs for the GP approximation leads to a so-called sparse spectrum GP, first proposed by Zato-Gredilla et al.

(2010) in the context of GP regression. Sparse spectrum GP with the offset term . More precisely, for SSPP, the hyperparameters consist of $(\gamma; \sigma^2)$, while for GSSPP, in Equations (10) or (13) $k^{(r)}(x; x^0) = \gamma^{(r)}(x) \gamma^{(r)}(x^0)$. $\gamma^{(r)} := (\gamma_{k=1}^{(r)}; \dots; \gamma_{k=r}^{(r)})$.

The resulting approximate Gaussian process can be written in terms of a new-size latent vector $w^{(r)}$ becomes

$$f^{(r)}(x) = w^{(r)\top} \gamma^{(r)}(x) \text{ where } w^{(r)} \sim N(0; I_r) \quad (14)$$

We define a permanent process with the spectral approximation of Equation (14), in which the intensity vector follows a similar linear form with feature vector given by Equations (10) or (13). We refer to it as Sparse Spectral Permanent Process (SSPP) when using the feature map (10) or Generalized Sparse Spectral Permanent Process (GSSPP) when using the feature map (13).

In both cases, a tractable expression can be obtained for the integral term in the likelihood, now defined as $\int_{\mathcal{X}} \gamma^{(r)}(x) dx := \int_{\mathcal{X}} (f^{(r)}(x) + \gamma^{(r)}(x))^2 dx$, as follows:

Proposition 3.3. (with proof in parts C and D of the Appendix) Under the GP approximation of Equation (14), the integral expression $\int_{\mathcal{X}} \gamma^{(r)}(x) dx$ can be expressed as

$$\int_{\mathcal{X}} \gamma^{(r)}(x) dx = w^{(r)\top} M^{(r)} w^{(r)} + 2 w^{(r)\top} m^{(r)} + \gamma^{(r)\top} X \gamma^{(r)} \quad (15)$$

where $M^{(r)}$ is an $r \times r$ matrix with i, j entries defined as

$$M_{ij}^{(r)} := \int_{\mathcal{X}^2} \gamma_i^{(r)}(x) \gamma_j^{(r)}(x) dx \quad (16)$$

for $i, j = 1; \dots; r$ and $m^{(r)}$ is a r -vector with entries

$$m_i^{(r)} := \int_{\mathcal{X}} \gamma_i^{(r)}(x) dx \quad (17)$$

for $i = 1; \dots; r$. Final expressions for $M^{(r)}$ and $m^{(r)}$ are provided in the Appendix

The solution of Proposition 3.3 shares similarities with Warren et al. (2022) approach, who utilized Random Fourier Features (RFF) in Bayesian Quadrature (BQ). However, our result also covers the calculation of $\int_{\mathcal{X}} \gamma^{(r)}(x) dx$ and applies to the feature map of a generalized kernel, with the feature map of standard RFF being a specific instance.

4 INFERENCE

Adopting the sparse spectral GP, $f^{(r)}$ assumes the linear form of Equation (14) for weight vector $w^{(r)}$ with independently-distributed standard Gaussian elements. Moreover, the integral in Equation (15) reduces to a quadratic form. We also define the model hyperparameters to be the parameters of the kernel function together

$$\log p(w^{(r)}; X, j) = \log p(X j w^{(r)};) + \log p(w^{(r)}) = \log p(X j w^{(r)};) - \frac{1}{2} w^{(r)\top} w^{(r)} + C \quad (18)$$

where $M^{(r)}$ and $m^{(r)}$ are the matrix and vector terms from Proposition 3.3. We can compute the log-likelihood function in $O(r^2 N)$, i.e. linearly in N . The log of the joint-distribution over $w^{(r)}$ and X is then

$$\log p(w^{(r)}; X, j) = \log p(X j w^{(r)};) + \log p(w^{(r)}) = \log p(X j w^{(r)};) - \frac{1}{2} w^{(r)\top} w^{(r)} + C \quad (19)$$

for some constant C , where $p(w^{(r)}) = N(0; I_r)$ denotes the prior distribution over $w^{(r)}$.

4.1 Laplace Approximation

The latent posterior $p(w^{(r)}; X, j)$ induced from Equation (19) is approximated using Laplace's method. A second order Taylor expansion of $\log p(w^{(r)}; X, j)$ around the maximum of the posterior, yields a Gaussian approximation

$$p(w^{(r)}; X, j) \approx N(w^{(r)}; \hat{w}^{(r)}, Q) := q(w^{(r)}; X, j) \quad (20)$$

where $\hat{w}^{(r)} := \arg \max_{w^{(r)}} p(w^{(r)}; X, j)$ is the mode of the latent posterior and Q is chosen to be the negative inverse Hessian of the true posterior at that point.

The gradient and the Hessian of the true posterior with respect to $w^{(r)}$ are

$$\begin{aligned} \nabla_{w^{(r)}} \log p(w^{(r)}; X, j) &= (2M^{(r)} + I_r) w^{(r)} + 2 m^{(r)} + 2 \sum_{i=1}^N \frac{\gamma^{(r)}(x_i)}{w^{(r)\top} \gamma^{(r)}(x_i) + \gamma^{(r)\top} X \gamma^{(r)}} \\ \nabla_{w^{(r)}}^2 \log p(w^{(r)}; X, j) &= (2M^{(r)} + I_r) - 2 \sum_{i=1}^N \frac{\gamma^{(r)}(x_i) \gamma^{(r)\top}(x_i)}{(w^{(r)\top} \gamma^{(r)}(x_i) + \gamma^{(r)\top} X \gamma^{(r)})^2} \end{aligned}$$

The mode $\hat{w}^{(r)}$ must satisfy the stationary constraint

$$\nabla_{w^{(r)}} \log p(w^{(r)}; X, j) \Big|_{w^{(r)} = \hat{w}^{(r)}} = 0;$$

(a) $\mu_1(x)$ (b) $\mu_2(x)$ (c) $\mu_3(x)$

Figure 2: Mean predictive intensity of the three toy intensity functions μ_2 and μ_3 defined as in Adams et al. (2009). Solid colored lines represent the predictive mean. The solid black lines shows the ground truth. The shaded areas are the 80% credible region of the SSPP model.

that implies

$$M^{(r)} + \frac{1}{2} \log \left(\prod_{i=1}^N \frac{\mu^{(r)}(x_i)}{w^{(r)}(x_i) + m^{(r)}} \right) \quad (21)$$

Equation (21) cannot be solved analytically. Instead, we estimate $w^{(r)}$ iteratively using Newton-Raphson method, with step

$$w^{(r)\text{new}} = w^{(r)} - \frac{r \frac{\partial}{\partial w^{(r)}} \log p(w^{(r)})}{r \frac{\partial}{\partial w^{(r)}} \log p(w^{(r)})} \quad (22)$$

The precision matrix is Q^{-1} is then given by $r \frac{\partial}{\partial w^{(r)}} \log p(w^{(r)} | X; \mu) |_{w^{(r)} = w^{(r)}}$.

4.2 Model Selection

We first derive a marginal likelihood approximation similar to Walder and Bishop (2017, Section 4.1.6).

$$\begin{aligned} \log p(X; \mu) &= \log p(w^{(r)}; X; \mu) - \log p(w^{(r)} | X; \mu) \\ &= \log p(w^{(r)}; X; \mu) - \log q(w^{(r)} | X; \mu) \\ &= \sum_{j=1}^N \log p(w^{(r)}; X_j; \mu) - \sum_{j=1}^N \log q(w^{(r)}; X_j; \mu) \\ &= \sum_{j=1}^N \left[\log p(w^{(r)}; X_j; \mu) - \log q(w^{(r)}; X_j; \mu) \right] \\ &= \sum_{j=1}^N \left[\log p(w^{(r)}; X_j; \mu) - \log q(w^{(r)}; X_j; \mu) \right] \\ &= \sum_{j=1}^N \left[\log p(w^{(r)}; X_j; \mu) - \log q(w^{(r)}; X_j; \mu) \right] \end{aligned} \quad (23)$$

since the quadratic term $\log q(w^{(r)} | X; \mu)$ cancels out.

We tune the hyperparameters by maximizing Equation (23). The model selection is facilitated by the fact that the gradient of the marginal likelihood in Equation (23) with respect to $w^{(r)}$ can be easily expressed. The terms $M^{(r)}$, $m^{(r)}$ and $\mu^{(r)}(\cdot)$ are functions of the hyperparameters. The mode $w^{(r)}$ is also a function of $\mu^{(r)}$.

The partial derivatives of the marginal likelihood with respect to $w^{(r)}$ is obtained using the chain rule,

$$r \frac{\partial}{\partial w^{(r)}} \log p(x_j) = \frac{\partial \log p(x_j)}{\partial \mu_j} \frac{\partial \mu_j}{\partial w^{(r)}} + \sum_{j=1}^N \frac{\partial \log p(X_j)}{\partial w_j} \frac{\partial w_j}{\partial w^{(r)}} \quad (24)$$

Expressions for the terms $\frac{\partial \log p(X_j)}{\partial \mu_j}$, $\frac{\partial \log p(X_j)}{\partial w_j}$ and $\frac{\partial \mu_j}{\partial w^{(r)}}$ above are given in part F of the Appendix, requiring a full mode search within each iterative hyperparameters update. In the current work, we note that assuming $\frac{\partial \mu_j}{\partial w^{(r)}} = 0$ and alternating independent updates for the mode in Equation (22) and the hyperparameters in Equation (24) provides faster and yet acceptable results. The algorithms are presented in part G of the Appendix.

5 PREDICTIVE DISTRIBUTION

To form predictive distributions, we assume that the latent posterior is approximated by $q(w^{(r)} | X; \mu)$ as in Equation (20).

5.1 Predictive Intensity Distribution

For some $x \in X$, the predictive distribution $f(x)$ can be deduced from Equations (14) and (20) to be

$$f(x) | X; \mu \sim N(\mu(x); \Sigma(x)) \quad (25)$$

where

$$\mu(x) := \mu^{(r)}(x) \quad (26)$$

and

$$\Sigma(x) := \mu^{(r)}(x)^T Q^{-1} \mu^{(r)}(x) \quad (27)$$

Given $\mu(x) = (\mu(x) + \sigma^2)^2$ and Equation (25), we can also derive the predictive distribution of the intensity function i.e.

$$f(x) | X; \mu \sim \text{Gamma}(a; b) \quad (28)$$

with parameters α and β expressed in part E of the Appendix.

5.2 Predictive Expected Log-likelihood

For a training set $X = \{x_i\}_{i=1}^N$ and an held-out test set $X_{test} = \{x_i\}_{i=1}^N$, we can derive from Equation (1), an approximation for the expected predictive log-likelihood

$$E[\log p(X_{test}|X)] = E_{w^{(r)}} \int \prod_{i=1}^N (w^{(r)}(x_i) + \beta)^{-1} dx + \sum_{i=1}^N E_{w^{(r)}} \log(w^{(r)}(x_i) + \beta)$$

where $w^{(r)}(x) = q(w^{(r)}(x)|X)$. The expectation over the integral term can be solved analytically. The sum-of-expectation can be expressed using Pochhammer series, that we approximate in practise by interpolation of a look-up table of precomputed values. This is very similar to Lloyd et al. (2015, section 4.3). We provide more details in part E of the Appendix.

6 EXPERIMENTS

We benchmark the SSPP scheme introduced in Sections 4 and 5 against a Nystrom-based implementation of the LBPP scheme of Walder and Bishop (2017), a frequentist kernel smoothing approach with edge correction (KS) (Diggle, 1985) and a variational inference scheme for point processes proposed by Lloyd et al. (2015), referred to as the Variational Bayesian Point Process (VBPP). We test the algorithms on three 1D synthetic data sets and three real data sets (one in 1D and two in 2D).

6.1 Benchmarks Settings

Our KS implementation uses standard kernel density estimation with truncated normal kernels to account for domain knowledge. The kernel bandwidth parameter is estimated via grid search using the leave-one-out log average likelihood objective of Lloyd et al. (2015). We used a publicly-available implementation of VBPP (<https://github.com/st-vbpps/vbpps>). We adopted fixed inducing points on a grid of size $\lfloor \sqrt{N} \rfloor$. For consistency, we also used a constant offset for both LBPP and VBPP implementations.

6.2 Performance Metrics

The average test expected log-likelihood $L_{test} := E[\log p(X_{test}|X)]$ is used as an evaluation metric. This is generally difficult to compute for point process models, but is available for SSPP and LBPP (see Section 5). For the synthetic experiment we also consider the normalized L_2 norm to the known ground truth intensity function i.e.

Table 1: Performance of GSSPP, SSPP, KS, VBPP and LBPP schemes on three samples of synthetic data. Values in bold-face refer to best performance, which corresponds to lower values of L_2 , but higher values of L_{test} .

	$1(x)$			$2(x)$			$3(x)$		
	L_2	L_{test}	time(s)	L_2	L_{test}	time(s)	L_2	L_{test}	time(s)
GSSPP-SE	0.74	105.92	1.32	0.83	52.62	1.59	1.68	841.77	2.32
GSSPP-m12	0.71	104.23	1.88	0.86	48.58	1.42	1.98	838.17	2.54
GSSPP-m32	0.68	106.05	1.79	0.78	51.65	1.56	1.75	840.80	2.53
GSSPP-m52	0.69	106.12	1.71	0.84	52.77	1.59	1.68	841.77	2.32
SSPP	0.78	105.19	0.63	0.74	56.01	0.69	1.60	835.35	1.05
KS	1.10	102.49	0.09	0.85	58.07	0.08	3.22	834.68	0.19
VBPP	0.72	104.65	1.98	0.85	51.11	1.15	1.63	838.65	1.52
LBPP	0.81	103.29	0.06	0.96	51.76	0.07	1.87	833.31	0.13

6.3 Synthetic Dataset

Three 1D simulated examples from Adams et al. (2009) are considered. The corresponding intensities are defined as $\lambda_1(x) = 2 \exp(-x/15) + \exp(-((x-25)/10)^2)$ on the interval $[0, 50]$ for approximately 47 events per sample, $\lambda_2(x) = 5 \sin(x^2) + 6$ on the interval $[0, 5]$ for approximately 36 events per sample and $\lambda_3(x)$ is a piecewise linear function shown in Figure 2 on the interval $[0, 100]$ for approximately 225 events per sample.

These intensity functions have been considered previously in the context of Gaussian Cox process by Samo and Roberts (2015b), Donner and Opper (2018), John and Hensman (2018) and Aglietti et al. (2019). We train the models on 10 independent samples generated from the ground truth, and evaluate the performance of each using 50 test sets sampled independently from the ground truth. We use the acronyms GSSPP-SE and GSSPP-m to refer to the generalized kernel variants of SSPP as in Equation (11) with κ to be the Gaussian kernel and Mat kernel with parameter ν respectively.

We report optimal performance across models for sets of spectral points or inducing points of size (denoted by ρ) ranging from 15 to 100. Results are given in Table 1 and the mean predictive intensities displayed in Figure 2. GSSPP and SSPP outperform the other methods in terms of both L_{test} and L_2 for λ_1 and λ_2 . Compared to LBPP, GSSPP and SSPP perform better consistently, but with slightly increase execution times. In a similar manner to findings for VBPP, GSSPP and SSPP fitting remains up to three orders of magnitude faster than alternative MCMC-based methods (Adams et al., 2009), see Lloyd et al. (2015) for comparison.

(a) GSSPP-m12 (b) SSPP (c) VBPP (d) LBPP

Figure 3: Heat map of the predictive mean intensity for Taxi data set scaled to a unit square. The black dots are the input data points.

the performances per number of spectral points or inducing points.

As a third dataset, we consider the Porto taxi dataset (Moreira-Matias et al., 2013) which contains $7 \cdot 10^6$ trajectories of taxi journeys in the years 2013-2014 in the Portuguese city of Porto. We consider the pick-up locations as observations of a point process. As in Aglietti et al. (2019), we restrict the analysis to 7000 events selected with (latitude, longitude) pairs bounded by $(41:147; 8:58)$ and $(41:18; 8:65)$. We select 400 events at random as training set and use the rest as testing set. We are ranging from 15 to 200 Table 2 presents the results. Figure 3 illustrates a single t to the full data set for four models. Figure 5 in Appendix H shows the performances per number of spectral points or inducing points.

Table 2: Results on real-world data experiments with standard errors in brackets.

Figure 4: Predictive mean intensity for the coal mine accident, with highest 80% credible intervals.

6.4 Real Datasets

The classic coal mine accidents data set consists of the dates of 191 coal-mining accidents with fatalities in Britain between 15 March 1875 and 22 March 1962 (Jarrett, 1979). For this data set, we evaluate predictive performance for the competing inference schemes using 100 random partitions of the sample into train and test subsets (and X) of approximately equal size. Figure 4 shows the predicted mean intensity with credible intervals. Results are presented in Table 2.

	Coal data (1D)		Bei data (2D)		Taxi data (2D)	
	L_{test}	time(s)	$L_{test} [10^{-1}]$	time(s)	$L_{test} [10^{-2}]$	time(s)
GSSPP-SE	224.44 (0.57)	1.56	763.49 (3.81)	20.48	278.54 (1.64)	23.91
GSSPP-m12	220.80 (0.85)	1.84	760.82 (4.31)	20.31	283.23 (1.11)	23.86
GSSPP-m32	224.25 (0.55)	1.58	764.73 (2.60)	20.85	280.48 (0:72)	19.60
GSSPP-m52	223.84 (0.54)	1.25	763.82 (1.00)	20.55	281.18 (1.43)	23.94
SSPP	221.23 (0.86)	0.64	751.50 (2.55)	17.54	268.32 (0.65)	12.13
KS	219.50 (0.33)	0.11	735.78 (1.49)	4.22	262.13 (0.26)	2.55
VBPP	221.19 (1.34)	1.75	757.95 (3.14)	28.41	281.02 (0.63)	20.45
LBPP	218.68 (0.87)	0.16	711.72 (1.35)	1.35	254.45 (0.17)	1.04

The bei data set is comprised of the locations of 3605 trees in the tropical rainforest on Barro Colorado Island (Hubbell and Foster, 1983). For these data, we evaluate predictive performance using 100 random partitions of the original sample into train and test subsamples of approximately equal size. Table 2 presents the results. Figure 6 in Appendix H, provides an illustration of larger coefficients for trigonometric components; hence the resulting intensity appears less smooth.

For each of the real-data applications, GSSPP performs best. The average fitting time using GSSPP is comparable to that of VBPP. Figure 4 shows the effect of different choice of GSSPP kernel function. When choosing α to be a Matern kernel, with parameter ν , spectral points x are drawn from a Student-t distribution with ν degree of freedom as discussed in part C of the Appendix. Compared to GSSPP with p now ranging from 15 to 150. Table 2 presents the results. Figure 6 in Appendix H, provides an illustration of larger coefficients for trigonometric components; hence the resulting intensity appears less smooth.

7 CONCLUSION

We introduce a novel Bayesian framework to infer the intensity function of a permanent process. Our approach uses a Laplace-based inference exploiting generalized kernels and random Fourier features (RFFs). The approach requires no discretization of the domain, allows kernel designs, and provides better predictive accuracy than the alternative Laplace-based approach of Walder and Bishop (2017). The performance of our scheme also compares favorably with other standard methods on both real temporal and large spatial data sets.

Acknowledgements

We also thank the anonymous reviewers for useful comments during the review process. Jeremy Sellier was supported under an EPSRC CASE studentship award in conjunction with Shell Research Ltd.

References

- Adams, R. P., Murray, I. and MacKay, D. J. (2009), Tractable nonparametric Bayesian inference in Poisson processes with Gaussian process intensities, *Proceedings of the 26th Annual International Conference on Machine Learning*, pp. 9–16.
- Aglietti, V., Bonilla, E. V., Damoulas, T. and Cripps, S. (2019), Structured variational inference in continuous Cox process models, *Advances in Neural Information Processing Systems*, Vol. 32.
- Banerjee, S., Carlin, B. P., Gelfand, A. E. and Banerjee, S. (2003), *Hierarchical Modeling and Analysis for Spatial Data*, Chapman and Hall/CRC. 1st edition.
- Bochner, S. (1932), *Vorlesungen über fouriersche integrale*, Akademische Verlagsgesellschaft.
- Cox, D. (1955), *Some statistical methods connected with series of events*, *Journal of the Royal Statistical Society Series B*, 17, 129–164.
- Cunningham, J. P., Shenoy, K. V. and Sahani, M. (2008), Fast Gaussian process methods for point process intensity estimation, in *Proceedings of the 25th International Conference on Machine Learning*, ICML '08, Association for Computing Machinery, p. 192–199.
- Cunningham, J. P., Yu, B. M., Shenoy, K. V. and Sahani, M. (2008), Inferring neural firing rates from spike trains using gaussian processes, *Advances in Neural Information Processing Systems*, Vol. 20.
- Daley, D. J. and Vere-Jones, D. (2003), *An introduction to the theory of point processes. Vol. I. Probability and its Applications*, Springer, New York.
- Diggle, P. J. (1985), *A Kernel method for smoothing point process data*, *Journal of The Royal Statistical Society Series C-applied Statistics*, 34, 138–147.
- Diggle, P. J. (2003), *Statistical analysis of spatial point patterns*, Hodder Education. 2nd edition.
- Diggle, P. J., Moraga, P., Rowlingson, B. and Taylor, B. M. (2013), *Spatial and spatio-temporal log-Gaussian Cox processes: extending the geostatistical paradigm*, *Statistical Science*, 28(4), 542 – 563.
- Donner, C. and Opper, M. (2018), *Efficient Bayesian inference of Sigmoidal Gaussian Cox processes*, *Learn. Res*, 19(1), 2710–2743.
- Duvenaud, D. K., Nickisch, H. and Rasmussen, C. (2011), *Additive Gaussian processes*, *Advances in Neural Information Processing Systems*, Vol. 24, Curran Associates, Inc.
- Duvenaud, D., Lloyd, J. R., Grosse, R., Tenenbaum, J. B. and Ghahramani, Z. (2013), Structure discovery in nonparametric regression through compositional kernel search, in *Proceedings of the 30th International Conference on International Conference on Machine Learning - Volume 28*, ICML'13.
- Flaxman, S., Chirico, M., Pereira, P. and Loeffler, C. (2019), *Scalable high-resolution forecasting of sparse spatiotemporal events with kernel methods: a winning solution to the NIJ “real-time crime forecasting challenge”*, *Annals of Applied Statistics*.
- Flaxman, S., Teh, Y. W. and Sejdinovic, D. (2017), *Poisson intensity estimation with reproducing kernel Hilbert spaces*, *Electronic Journal of Statistics*, 1(2), 5081–5104.
- Flaxman, S., Wilson, A., Neill, D., Nickisch, H. and Smola, A. (2015), *Fast Kronecker inference in Gaussian processes with non-Gaussian likelihoods*, *Proceedings of the 32nd International Conference on Machine Learning*, Vol. 37 of *Proceedings of Machine Learning Research*, pp. 607–616.
- Gardner, J. K. and Knopoff, L. (1974), *Is the sequence of earthquakes in Southern California, with aftershocks removed, Poissonian?*, *Bulletin of the Seismological Society of America*, 64(5), 1363–1367.
- Grubisic, T. H. and Mack, E. A. (2008), *Spatio-temporal interaction of urban crime*, *Journal of Quantitative Criminology*, 24, 285–306.
- Gunter, T., Lloyd, C., Osborne, M. and Roberts, S. (2014), *Efficient Bayesian nonparametric modelling of structured point processes*, *Uncertainty in Artificial Intelligence - Proceedings of the 30th Conference, UAI 2014*.
- Hu, S. and Papageorgiou, N. (2018), *Handbook of Multivalued Analysis: Volume I: Theory*, Mathematics and Its Applications, Springer US.
- Hubbell, S. and Foster, R. (1983), *Diversity of canopy trees in a neotropical forest and implications for conservation*, *Tropical Rain Forest: Ecology and Management*, 11, 25–41.

- Illian, J. B., Sørbye, S. H. and Rue, H. (2012), 'A toolbox for fitting complex spatial point process models using integrated nested Laplace approximation (INLA)', *The Annals of Applied Statistics*, 4(4), 1499–1530.
- Jarrett, R. G. (1979), 'A note on the intervals between coal-mining disasters', *Biometrika*, 66(1), 191–193.
- John, S. and Hensman, J. (2018), 'Large-scale Cox process inference using variational Fourier features', *International Conference on Machine Learning*, pp. 2362–2370.
- Ko, Y.-J. and Seeger, M. W. (2016), 'Expectation propagation for rectified linear poisson regression', *Asian Conference on Machine Learning*, Vol. 45 of *Proceedings of Machine Learning Research*, Hong Kong, pp. 253–268.
- Lázaro-Gredilla, M., Quinonero-Candela, J., Rasmussen, C. E. and Figueiras-Vidal, A. R. (2010), 'Sparse spectrum Gaussian process regression', *The Journal of Machine Learning Research*, 11, 1865–1881.
- Lian, W., Henao, R., Rao, V., Lucas, J. p. and Carin, L. (2015), 'A multitask point process predictive model', *Proceedings of the 32nd International Conference on Machine Learning*, Vol. 37 of *CML'15*, pp. 2030–2038.
- Lloyd, C., Gunter, T., Osborne, M. and Roberts, S. (2015), 'Variational inference for Gaussian process modulated Poisson processes', *International Conference on Machine Learning*, pp. 1814–1822.
- Lopez-lopera, A. F., John, S. and Durrande, N. (2019), 'Gaussian process modulated cox processes under linear inequality constraints', *Proceedings of the Twenty-Second International Conference on Artificial Intelligence and Statistics*, Vol. 89 of *Proceedings of Machine Learning Research*, pp. 1997–2006.
- McCullagh, P. and Møller, J. (2006), 'The permanental process', *Advances in Applied Probability*, 38(4), 873–888.
- Moreira-Matias, L., Gama, J., Ferreira, M., Mendes-Moreira, J. and Damas, L. (2013), 'Predicting taxi-passenger demand using streaming data', *IEEE Transactions on Intelligent Transportation Systems*, 14(3), 1393–1402.
- Møller, J., Syversveen, A. R. and Waagepetersen, R. P. (1998), 'Log Gaussian Cox processes', *Scandinavian Journal of Statistics*, 25(3), 451–482.
- Park, M., Weller, J., Horwitz, G. and Pillow, J. (2014), 'Bayesian active learning of neural firing rate maps with transformed gaussian process priors', *Neural computation*, 26, 1–23.
- Plataniotis, K. N. and Hatzinakos, D. (2003), *Gaussian Mixtures and their Applications to Signal Processing*. Advanced Signal Processing Handbook, CRC Press Springer-Verlag.
- Rahimi, A. and Recht, B. (2007), 'Random features for large-scale kernel machines', *Proceedings of the 20th International Conference on Neural Information Processing Systems*, p. 1177–1184.
- Rasmussen, C. E. and Williams, C. K. I. (2006), *Gaussian processes for machine learning*. MIT Press.
- Samo, Y.-L. K. (2017), 'Advances in kernel methods : towards general-purpose and scalable models', PhD thesis, University of Oxford.
- Samo, Y.-L. K. and Roberts, S. (2015), 'Generalized spectral kernels', arxiv:1506.02236
- Samo, Y.-L. K. and Roberts, S. (2015), 'Scalable nonparametric Bayesian inference on point processes with Gaussian processes', *Proceedings of the 32th International Conference on Machine Learning*, Vol. 37 of *CML'15*, pp. 2227–2236.
- Seeger, M. and Bouchard, G. (2012), 'Fast variational bayesian inference for non-conjugate matrix factorization models', in *Proceedings of the 15th International Conference on Artificial Intelligence and Statistics*, Vol. 22 of *Proceedings of Machine Learning Research*, pp. 1012–1018.
- Stein, M. L. (1999), *Interpolation of spatial data*. Springer Series in Statistics, Springer-Verlag, New York.
- Titsias, M. (2009), 'Variational learning of inducing variables in sparse Gaussian processes', *Proceedings of the 20th International Conference on Artificial Intelligence and Statistics*, Vol. 5 of *Proceedings of Machine Learning Research*, pp. 567–574.
- Walder, C. J. and Bishop, A. N. (2017), 'Fast Bayesian intensity estimation for the permanental process', *Proceedings of the 34th International Conference on Machine Learning*, Vol. 70 of *Proceedings of Machine Learning Research*, pp. 3579–3588.
- Warren, H., Oliveira, R. and Ramos, F. (2022), 'Generalized bayesian quadrature with spectral kernels', *Proceedings of the 38th Conference on Uncertainty in Artificial Intelligence*, Vol. 180 of *Proceedings of Machine Learning Research*, pp. 2085–2095.
- Wilson, A. and Adams, R. (2013), 'Gaussian process kernels for pattern discovery and extrapolation', *Proceedings of the 30th International Conference on Machine Learning*, pp. 1067–1075.

A LBPP WITH NYSTRÖM METHOD

We provide a quick review of the LBPP model using Nyström method proposed by Walder and Bishop (2017). We also provide in Proposition A.1, an expression for the integral term $\int_{\mathcal{X}} f(x)^2 dx$ under Nyström approximation that is not directly available in Walder and Bishop (2017).

A.1 Method Review

The Nyström method (see Rasmussen and Williams, 2005, Chap. 4.3, 8) provides a reduced-rank approximation for $n \ll N$ data points $X^{(n)} = \{x_i^{(n)}\}_{i=1}^n$ sampled uniformly from the original data. The eigenvalues and eigenfunctions in Equation (3) of the main text are approximated using the eigenvectors and eigenvalues $\lambda_i^{(n)}$ of $K_{n,n}$, the Gram matrix with i, j entry $k(x_i^{(n)}; x_j^{(n)})$. Thus, f can be approximated by

$$f(x) \approx \sum_{i=1}^n \lambda_i^{-\frac{1}{2}} w_i^{(n)} \hat{\phi}_i(x) \quad (29)$$

where $w^{(n)} = (w_1^{(n)}; \dots; w_n^{(n)})^T \in \mathbb{R}^n$ and

$$\hat{\phi}_i := \frac{1}{\sqrt{\lambda_i^{(n)}}} \sum_{j=1}^n \frac{1}{N} k(x_j^{(n)}; x_i^{(n)}) u_j^{(n)} \quad (30)$$

$$u_i^{(n)} := \frac{1}{\sqrt{\lambda_i^{(n)}}} \sum_{j=1}^n k(x_j^{(n)}; x_i^{(n)}) u_j^{(n)} \quad (31)$$

Substituting Equations (30) and (31) into Equation (29) yields a linear Gaussian latent formulation as follows

$$\begin{aligned} f(x) &\approx k(x; X^{(n)}) \sum_{i=1}^n \frac{w_i^{(n)}}{\sqrt{\lambda_i^{(n)}}} u_i^{(n)} \\ &= k(x; X^{(n)}) U^{(n)} \left(\sum_{i=1}^n \frac{1}{\sqrt{\lambda_i^{(n)}}} w_i^{(n)} \right) \\ &:= w^{(n)T} \phi^{(n)}(x) \end{aligned} \quad (32)$$

and $\phi^{(n)}(x) := \frac{1}{\sqrt{\lambda_i^{(n)}}} \sum_{i=1}^n k(x; X^{(n)}) U^{(n)} u_i^{(n)}$ denotes the new features vector.

A.2 Integral Calculation

In Proposition A.1 we express the integral term $\int_{\mathcal{X}} f(x)^2 dx$ under the Nyström approximation both because it is not available in Walder and Bishop (2017) and to demonstrate the similarities with the corresponding derivation of our proposed method in Proposition 3.3.

Proposition A.1. Under the GP approximation (32), the integral expression $\int_{\mathcal{X}} f(x)^2 dx$ can be written as

$$\int_{\mathcal{X}} f(x)^2 dx = w^{(n)T} M^{(n)} w^{(n)}$$

where $M^{(n)}$ is an $n \times n$ matrix defined as

$$M^{(n)} := \left(\sum_{i=1}^n \frac{1}{\lambda_i^{(n)}} U^{(n)} u_i^{(n)} u_i^{(n)T} \right)$$

and $\phi^{(n)}$ is an n matrix given in the proof below.

Proof. Let f be approximated by the Nyström-based approach defined in Equation (32) $f(x) \approx w^{(n)T} \phi^{(n)}(x)$. The

integral expression $\int_{\mathcal{X}} f(x)^2 dx$ can be written as

$$\begin{aligned} \int_{\mathcal{X}} f(x)^2 dx &= \int_{\mathcal{X}} \sum_{i=1}^n \sum_{j=1}^n w_i^{(n)} w_j^{(n)} \phi_i^{(n)}(x) \phi_j^{(n)}(x) dx \\ &= \sum_{i=1}^n \sum_{j=1}^n \frac{w_i^{(n)} w_j^{(n)}}{u_i^{(n)} u_j^{(n)}} \int_{\mathcal{X}} k(x; X^{(n)}) u_i^{(n)} k(x; X^{(n)}) u_j^{(n)} dx \\ &= \sum_{i=1}^n \sum_{j=1}^n \frac{w_i^{(n)} w_j^{(n)}}{u_i^{(n)} u_j^{(n)}} u_i^{(n)} u_j^{(n)} \int_{\mathcal{X}} k(X^{(n)}; x) k(x; X^{(n)}) dx \\ &= \mathbf{w}^{(n)\top} \mathbf{U}^{(n)} \mathbf{U}^{(n)\top} \mathbf{w}^{(n)} \end{aligned}$$

where $u_i^{(n)} = \int_{\mathcal{X}} k(X^{(n)}; x) k(x; X^{(n)}) dx$ is the integral statistic already defined in Lloyd et al. (2015) and John and Hensman (2018). In particular, for the separable Gaussian kernel defined in Equation (35),

$$\begin{aligned} u_{ij} &= \int_{\mathcal{X}} \prod_{k=1}^d \exp\left(-\frac{(x_{k,i}^{(n)} - x_{k,j}^{(n)})^2}{4\sigma_k^2}\right) \exp\left(-\frac{(x_{k,i}^{(n)} - x_{k,j}^{(n)})^2}{\sigma_k^2}\right) dx \\ &= \prod_{k=1}^d \frac{\sigma_k^2}{d} \exp\left(-\frac{(x_{k,i}^{(n)} - x_{k,j}^{(n)})^2}{4\sigma_k^2}\right) \operatorname{erf}\left(\frac{x_{k,i}^{(n)} - x_{k,j}^{(n)}}{\sigma_k}\right) \operatorname{erf}\left(\frac{x_{k,i}^{(n)} - x_{k,j}^{(n)}}{\sigma_k}\right) \end{aligned}$$

where σ_k and $\sigma_k' := (\sigma_k^2, \sigma_k'^2)$ are respectively the scaling and length-scale parameters of the covariance function in the k th coordinate of the i th Nystöm-sampled point $\mathbf{x}_i^{(n)}$ and $\mathbf{x}_{k,i,j}^{(n)} := (x_{k,i}^{(n)} + x_{k,j}^{(n)})/2$.

Offset Term Adding an offset term to the intensity i.e. $f(x) = (f(x) + \beta)^2$ yields

$$\int_{\mathcal{X}} (f(x) + \beta)^2 dx = \int_{\mathcal{X}} f(x)^2 dx + 2 \int_{\mathcal{X}} f(x) dx + \beta^2 |\mathcal{X}|$$

with

$$\int_{\mathcal{X}} f(x) dx = \sum_{i=1}^n \sum_{j=1}^n \frac{w_i^{(n)} w_j^{(n)}}{u_i^{(n)} u_j^{(n)}} \int_{\mathcal{X}} k(x; X^{(n)}) u_i^{(n)} k(x; X^{(n)}) u_j^{(n)} dx = \mathbf{w}^{(n)\top} \mathbf{U}^{(n)} \mathbf{U}^{(n)\top} \mathbf{w}^{(n)}$$

where $\mathbf{w}^{(n)} := \sum_{i=1}^n \sum_{j=1}^n \frac{w_i^{(n)} w_j^{(n)}}{u_i^{(n)} u_j^{(n)}} \int_{\mathcal{X}} k(x; X^{(n)}) u_i^{(n)} k(x; X^{(n)}) u_j^{(n)} dx$ is an n -vector such that in the separable Gaussian kernel case above, we have

$$w_i^{(n)} = \prod_{d=1}^d \frac{\sigma_d^2}{d} \operatorname{erf}\left(\frac{x_{d,i}^{(n)} - x_{d,i}^{(n)}}{\sigma_d}\right) \operatorname{erf}\left(\frac{x_{d,i}^{(n)} - x_{d,i}^{(n)}}{\sigma_d}\right)$$

□

A.3 Inference with Laplace Approximation

For a fixed set of GP hyperparameters (i.e. parameters of the covariance function), the intensity function is determined by the latent vector $\mathbf{w}^{(n)}$. Walder and Bishop (2017) use a Laplace approximation for a non-Gaussian posterior $p(\mathbf{w}^{(n)} | \mathbf{X}; \theta)$. The latent posterior is approximated by $\mathbf{w}^{(n)} | \mathbf{X}; \theta \approx \mathcal{N}(\mathbf{w}^{(n)} | \mathbf{w}^{(n)}, \mathbf{Q})$ where $\mathbf{w}^{(n)}$ is the mode of the posterior and \mathbf{Q} is chosen to be the inverse Hessian of the true posterior at $\mathbf{w}^{(n)}$.

B COMMON KERNELS SPECTRAL DENSITIES

We present in Table (3) below some common stationary distance-dependent kernels mentioned in the main text and their corresponding spectral densities $S(\cdot)$.

Table 3: Stationary distance-dependent kernels and their duals, where $\Gamma(\cdot)$ is the Gamma function and $J_\nu(\cdot)$ is a modified Bessel function.

NAME	KERNEL FUNCTION $k(\cdot)$	SPECTRAL DENSITY $S(z)$
Gaussian	$\exp\left(-\frac{\ z\ ^2}{2}\right)$	$\frac{1}{(2\pi)^{d/2}} \exp\left(-\frac{\ z\ ^2}{2}\right)$
Matérn (ν)	$\frac{2^{1-\nu}}{\Gamma(\nu)} \left(\frac{\ z\ }{2}\right)^{\nu-\frac{d}{2}} K_{\nu-\frac{d}{2}}\left(\frac{\ z\ }{2}\right)$	$\frac{(\nu + d/2)^{-\nu}}{\Gamma(\nu)(2\pi)^{d/2}} \left(1 + \frac{\ z\ }{2}\right)^{\nu-\frac{d}{2}}$

C PROOF OF PROPOSITION (3.3) : Integral Expression via RFF

Let f be approximated by a RFF-based approach as defined in Equation (14) $f(x) \approx w^{(r)\top} \phi^{(r)}(x)$ where the feature map $\phi^{(r)}$ follows Equation (9).

C.1 Real Valued Feature Mapping

We first detail the derivation of the real valued Fourier features described in Equation (10) of the main text. The imaginary part of Equation (9) of the main text can be discarded as follows

$$\begin{aligned} k(x, x^0) &= \mathbb{E}_z \exp(i z^\top (x - x^0)) \\ &= \mathbb{E}_z \cos(z^\top (x - x^0)) + i \sin(z^\top (x - x^0)) \end{aligned} \tag{33}$$

$$\begin{aligned} &= \mathbb{E}_z \cos(z^\top (x - x^0)) \\ &= \mathbb{E}_z \cos(z^\top x) \cos(z^\top x^0) + \sin(z^\top x) \sin(z^\top x^0) \end{aligned} \tag{34}$$

$$\begin{aligned} &= \frac{1}{r} \sum_{i=1}^r \cos(z_i^\top x) \cos(z_i^\top x^0) + \sin(z_i^\top x) \sin(z_i^\top x^0) \\ &= \frac{1}{r} \phi^{(r)\top}(x) \phi^{(r)}(x^0) \end{aligned}$$

where z_1, \dots, z_r are independent samples with density $S(z)$ and the explicit feature mapping $\phi^{(r)}(\cdot)$ is defined as

$$\phi^{(r)}(x) := \frac{1}{r} \begin{bmatrix} \cos(z_1^\top x) \\ \vdots \\ \cos(z_r^\top x) \\ \sin(z_1^\top x) \\ \vdots \\ \sin(z_r^\top x) \end{bmatrix}$$

Gaussian Kernel Specifically, without approximation, for a Gaussian kernel with $X = \mathbb{R}^d$ and where

$$k_g(x, x^0) = \frac{1}{(2\pi)^{d/2}} \exp\left(-\frac{\sum_{i=1}^d (x_i - x_i^0)^2}{2\sigma_i^2}\right) \tag{35}$$

with scaling parameter and length-scale vector $\sigma = [\sigma_1, \dots, \sigma_d]^\top$, the corresponding spectral density $S(z)$ is a multivariate normal $\mathcal{N}(0; \text{diag}(\sigma^2))$ with $\sigma^2 := [1/\sigma_1^2, \dots, 1/\sigma_d^2]^\top$.

Matérn Kernel For a Matérn class of kernel function k_m such that

$$k_m(x, x^0) = \frac{2^{2\nu}}{\Gamma(\nu)} \frac{K_\nu(\sqrt{2\nu} \sqrt{\lambda} \|x - x^0\|)}{(\sqrt{2\nu} \sqrt{\lambda} \|x - x^0\|)^\nu} \quad (36)$$

where $\nu \in \mathbb{R}^+$, $\lambda \in \mathbb{R}^+$ and K_ν is a modified Bessel function, the corresponding spectral density is a dimension multivariate Student-t distribution $\mathcal{N}(0; \Sigma^{-1})$ with covariance function $\Sigma = (1 + \lambda \|x - x^0\|)^{-2\nu} I_d$ and degree of freedom 2ν . The spectral locations \mathbf{z} are sampled as

$$\mathbf{Z} = \frac{1}{\sqrt{2\nu}} \mathbf{G} \quad \text{where } \mathbf{u} \sim \mathcal{N}(0, \mathbf{I}_d) \quad (37)$$

and \mathbf{G} is a $d \times r$ matrix of i.i.d. standard normal random variables.

C.2 Integral Calculation

We now detail the integral expression of proposition (3.3) for the real valued Fourier features in Equation (10) of the main text. We consider without loss of generality the spatial case where $\mathbf{x}_{d,i}$ refers to the i th coordinate of the training input \mathbf{x}_i for $i = 1, \dots, N$ and $z_{d,i}$ to the i th coordinate of the i th spectral point \mathbf{z}_i for $i = 1, \dots, r$. The integral of f over \mathbf{X} becomes

$$\int_{[a,a]^2} f(\mathbf{x})^2 dx = \sum_{i,j} w_i^{(r)} w_j^{(r)} \int_{[a,a]^2} \phi_i^{(r)}(\mathbf{x}) \phi_j^{(r)}(\mathbf{x}) dx \quad (38)$$

where

$$\phi_i^{(r)}(\mathbf{x}) \phi_j^{(r)}(\mathbf{x}) = \frac{1}{r} \begin{cases} \cos(z_i^\top \mathbf{x}) \cos(z_j^\top \mathbf{x}) & \text{if } (i,j) \in [1;r]^2 \\ \sin(z_i^\top \mathbf{x}) \sin(z_j^\top \mathbf{x}) & \text{if } (i,j) \in [r+2;2r]^2 \\ \cos(z_i^\top \mathbf{x}) \sin(z_j^\top \mathbf{x}) & \text{if } (i,j) \in [1;r] \times [r+2;2r] \\ \sin(z_i^\top \mathbf{x}) \cos(z_j^\top \mathbf{x}) & \text{if } (i,j) \in [r+2;2r] \times [1;r] \end{cases} \quad (39)$$

Thus, $\int_{\mathbf{X}} f(\mathbf{x})^2 dx = \mathbf{w}^\top \mathbf{M}^{(r)} \mathbf{w}$, where $\mathbf{M}^{(r)}$ is the matrix with (i,j) entry obtained by integrating Equation (39). The 'cos', 'sin' and 'cos-sin' expressions can be written as

$$\int_{[a,a]^2} \cos(z_i^\top \mathbf{x}) \cos(z_j^\top \mathbf{x}) dx = \frac{1}{2} \int_{[a,a]^2} \cos((z_i - z_j)^\top \mathbf{x}) + \cos((z_i + z_j)^\top \mathbf{x}) dx$$

$$\int_{[a,a]^2} \sin(z_i^\top \mathbf{x}) \sin(z_j^\top \mathbf{x}) dx = \frac{1}{2} \int_{[a,a]^2} \cos((z_i - z_j)^\top \mathbf{x}) - \cos((z_i + z_j)^\top \mathbf{x}) dx$$

and

$$\int_{[a,a]^2} \cos(z_i^\top \mathbf{x}) \sin(z_j^\top \mathbf{x}) dx = \frac{1}{2} \int_{[a,a]^2} \sin((z_i - z_j)^\top \mathbf{x}) + \sin((z_i + z_j)^\top \mathbf{x}) dx = 0:$$

Thus, since the off-diagonal blocks $\mathbf{M}^{(r)}$ are null, we can rewrite Equation (38) as

$$\int_{\mathbf{X}} f(\mathbf{x})^2 dx = \frac{1}{r} \mathbf{w}_{:r}^{(r)\top} (\mathbf{A} + \mathbf{B}) \mathbf{w}_{:r}^{(r)} + \mathbf{w}_{:r}^{(r)\top} (\mathbf{A} - \mathbf{B}) \mathbf{w}_{:r}^{(r)}$$

$$= \frac{1}{r} \mathbf{w}^{(r)\top} (\mathbf{D}_1^\top (\mathbf{A} + \mathbf{B}) \mathbf{D}_1 + \mathbf{D}_r^\top (\mathbf{A} - \mathbf{B}) \mathbf{D}_r) \mathbf{w}^{(r)} \quad (40)$$

where $\mathbf{w}_{:r}^{(r)} := [\mathbf{w}_1^{(r)}; \dots; \mathbf{w}_r^{(r)}]^\top$, $\mathbf{w}_{:r}^{(r)} := [\mathbf{w}_{r+1}^{(r)}; \dots; \mathbf{w}_{2r}^{(r)}]^\top$, $\mathbf{D}_1 := \begin{bmatrix} \mathbf{I}_r & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix}$, $\mathbf{D}_r := \begin{bmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_r \end{bmatrix}$ and \mathbf{A} and \mathbf{B} are two $r \times r$ matrices defined as

$$\mathbf{A}_{ij} = \frac{1}{2} \int_{[a,a]^2} \cos((z_i - z_j)^\top \mathbf{x}) dx \quad \text{and} \quad \mathbf{B}_{ij} = \frac{1}{2} \int_{[a,a]^2} \cos((z_i + z_j)^\top \mathbf{x}) dx:$$

\mathbf{A} and \mathbf{B} can be evaluated as follows

Case 1: $z_i \notin z_j$ We define again $z_{d;ij} := z_{d;i} + z_{d;j}$. Then,

$$A_{ij} = \frac{\cos[a(z_{1;ij} - z_{2;ij})] \cos[a(z_{1;ij} + z_{2;ij})]}{z_{1;ij} z_{2;ij}} = \frac{2}{z_{1;ij} z_{2;ij}} (\sin[az_{1;ij}] \sin[az_{2;ij}]) :$$

$$B_{ij} = \frac{\cos[a(z_{1;ij} - z_{2;ij})] \cos[a(z_{1;ij} + z_{2;ij})]}{z_{1;ij} z_{2;ij}} = \frac{2}{z_{1;ij} z_{2;ij}} (\sin[az_{1;ij}] \sin[az_{2;ij}]) :$$

Case 2: $z_i = z_j$

$$A_{ij} = 2a^2 :$$

$$B_{ij} = \frac{\cos[2a(z_{1;i} - z_{2;i})] \cos[2a(z_{1;i} + z_{2;i})]}{4z_{1;i} z_{2;i}} = \frac{1}{2z_{1;i} z_{2;i}} (\sin[2az_{1;i}] \sin[2az_{2;i}]) :$$

Offset Term For the offset term, we need to compute the integral that is obtained from

$$\int_{[a;a]^2} f(x) dx = \prod_{r=1}^R w_i^{(r)} \int_{[a;a]^2} \cos(z_i^> x) dx + \prod_{r=r+1}^R w_i^{(r)} \int_{[a;a]^2} \sin(z_i^> x) dx$$

$$= \prod_{r=1}^R w_i^{(r)} m^{(r)}$$

where m is an r -vector such that

$$m_i^{(r)} = \frac{2 \cos[a(z_{1;i} - z_{2;i})] \cos[a(z_{1;i} + z_{2;i})]}{z_{1;i} z_{2;i}} = \frac{4}{z_{1;i} z_{2;i}} (\sin[az_{1;i}] \sin[az_{2;i}]) :$$

D PROOF OF PROPOSITION (3.3) : Integral Expression for Generalized Kernel

We assume k_{GS} to be a Generalized kernel given in Equation (11) in the main text, with a kernel that admits a consistent RFF representation such that $k(x, x^0) = \prod_{g=1}^G \langle g^{(r)}(x) \rangle \langle g^{(r)}(x^0) \rangle$ where $\langle g^{(r)} \rangle$ is an explicit feature mapping $g : X \rightarrow \mathbb{R}^r$.

D.1 Real Valued Feature Mapping

The Generalized kernel k_{GS} becomes

$$k_{GS}(x, x^0) = \prod_{k=1}^K \langle k^{(r)}(x) \rangle \langle k^{(r)}(x^0) \rangle = \prod_{k=1}^K \langle k(x) \rangle \langle k(x^0) \rangle$$

where $\langle k(x) \rangle$ is a map $\langle k : X \rightarrow \mathbb{R}^2$ such that $\langle k(x) \rangle = \frac{\cos(\langle k, x \rangle)}{\sin(\langle k, x \rangle)}$ for $k = 1; \dots; K$ and $\forall x \in X$ so that $\langle k(x) \rangle \langle k(x^0) \rangle = \cos(\langle x - x^0, k \rangle)$ for all $k \in \mathbb{R}^d$ and $\forall x, x^0 \in X$. Thus,

$$k_{GS}(x; x^0) = \prod_{k=1}^K h_k(x) \langle h_k(x^0) \rangle$$

with

$$h_k(x) = \langle k, g^{(r)}(x) \rangle \frac{\cos(\langle k, x \rangle)}{\sin(\langle k, x \rangle)}$$

for $k = 1; \dots; K$, where $\langle \cdot, \cdot \rangle$ denotes the Kronecker product.

In particular, when $\langle g^{(r)} \rangle$ follows Equation (10) in the main text,

$$h_k(x) = \prod_{r=1}^R \left(\cos(z_1^>(x, k)) \frac{\cos(\langle k, x \rangle)}{\sin(\langle k, x \rangle)} \right) \prod_{r=1}^R \left(\sin(z_r^>(x, k)) \frac{\cos(\langle k, x \rangle)}{\sin(\langle k, x \rangle)} \right) \tag{41}$$

where z_1, \dots, z_r are independent samples from $\mathcal{G}(z)$ the spectral density of

The resulting approximate Gaussian process with generalized kernel can be written, in terms of a $4Kr$ size latent vector as follow

$$f(x) = w^{(r)\top} \mu^{(r)}(x) \quad \text{with } w^{(r)} \sim \mathcal{N}(0; I_{4Kr}):$$

where

$$\mu^{(r)}(x) = [h_1(x)^\top; \dots; h_K(x)^\top]^\top$$

with $h_k(x)$ defined as in Equation (41) for $k = 1; \dots; K$.

D.2 Integral Calculation

The integral over X becomes

$$\int_{[a,a]^2} f(x)^2 dx = \sum_{i,j} w_i^{(r)} w_j^{(r)} \int_{[a,a]^2} \mu_i^{(r)}(x) \mu_j^{(r)}(x) dx \quad (42)$$

Thus, $\int_{[a,a]^2} f(x)^2 dx = w^{(r)\top} M^{(r)} w^{(r)}$, where $M^{(r)}$ is the matrix with i,j entry obtained by integrating Equation (42). The computation of $M^{(r)}$ can be split into different cases expressed below as 'cos', 'sin' and 'cos-sin' terms.

Cos Terms The 'cos' i,j terms can be written as

$$\begin{aligned} & \int_{[a,a]^2} \cos(z_i^\top(x - z_i)) \cos(x^\top z_i) \cos(z_j^\top(x - z_j)) \cos(x^\top z_j) dx \\ &= \frac{1}{8} \sum_{k=1}^8 \int_{[a,a]^2} \cos(x^\top z_{ij}^{(k)}) dx \end{aligned} \quad (43)$$

where

$$\begin{aligned} (1)_{ij} &= (z_i - z_j) + \|z_i\| + \|z_j\|; & (5)_{ij} &= (z_i - z_j) - \|z_i\| + \|z_j\|; \\ (2)_{ij} &= (z_i - z_j) + \|z_i\| - \|z_j\|; & (6)_{ij} &= (z_i - z_j) - \|z_i\| + \|z_j\|; \\ (3)_{ij} &= (z_i - z_j) + \|z_i\| - \|z_j\|; & (7)_{ij} &= (z_i - z_j) - \|z_i\| + \|z_j\|; \\ (4)_{ij} &= (z_i - z_j) + \|z_i\| - \|z_j\|; & (8)_{ij} &= (z_i - z_j) - \|z_i\| + \|z_j\|. \end{aligned}$$

Integrating the left hand integrants in Equation (43) yields

$$\int_{[a,a]^2} \cos(x^\top z_{ij}^{(k)}) dx = \begin{cases} 8 < 2a^2; & \text{if } i = j \text{ and } k \in \{3, 7\}; \\ \frac{1}{2} \frac{\sin[a z_{1,ij}^{(k)}]}{z_{1,ij}^{(k)}} \sin[a z_{2,ij}^{(k)}] & \text{otherwise} \end{cases} \quad (44)$$

Sin Terms The 'sin' i,j terms are

$$\begin{aligned} & \int_{[a,a]^2} \sin(x^\top(z_i - z_j)) \sin(x^\top z_i) \sin(x^\top(z_j - z_i)) \sin(x^\top z_j) dx \\ &= \frac{1}{8} \sum_{k=1}^8 (-1)^k \int_{[a,a]^2} \cos(x^\top z_{ij}^{(k)}) dx \end{aligned} \quad (45)$$

The left hand integrants in Equation (45) integrate alike to Equation (44) up to a factor.

Cos-sin Terms The cos-sin terms can be evaluated as follows

$$\int_{[a;a]^2} \sin(x^> (z_i \quad i)) \sin(x^> ! i)) \cos(x^> (z_j \quad j)) \cos(x^> ! j)) dx$$

$$= \frac{1}{8} \sum_{k=1}^8 (1)^{m(k)} \int_{[a;a]^2} \cos(x^> \binom{k}{ij}) dx \quad (46)$$

where $m(k) = 1$ if $k = 1; 4$ and 0 else. The left hand integrands in Equation (46) integrate alike to Equation (44) up to a $m(k)$ factor.

The remaining terms, yields sums of integrals of the type $\int_{[a;a]^2} \sin(x^>) dx$ with $\int_{\mathbb{R}^d}$, that equal zero.

Offset Term For the offset term, we need to compute the integral of that is $\int_{\mathbb{R}^d} f(x) dx = w^{(r)} m^{(r)}$ where $m^{(r)}$ is a $4K^r$ -vector such that

$$m_i^{(r)} = \int_{[a;a]^2} f_i^{(r)}(x) dx$$

The computation of $m^{(r)}$ can be split into two cases : the cos' terms

$$\int_{[a;a]^2} \cos(z_i^> (x \quad i)) \cos(x^> ! i)) dx =$$

$$\frac{1}{\binom{(1)}{1;i} \binom{(1)}{2;i}} \sin[a \binom{(1)}{1;i}] \sin[a \binom{(1)}{1;i}] + \frac{1}{\binom{(2)}{1;i} \binom{(2)}{2;i}} \sin[a \binom{(2)}{1;i}] \sin[a \binom{(2)}{1;i}]$$

and the sin' terms

$$\int_{[a;a]^2} \sin(z_i^> (x \quad i)) \sin(x^> ! i)) dx =$$

$$\frac{1}{\binom{(1)}{1;i} \binom{(1)}{2;i}} \sin[a \binom{(1)}{1;i}] \sin[a \binom{(1)}{1;i}] - \frac{1}{\binom{(2)}{1;i} \binom{(2)}{2;i}} \sin[a \binom{(2)}{1;i}] \sin[a \binom{(2)}{1;i}]$$

where

$$\binom{(1)}{ij} = (z_i \quad i) + ! i; \quad \binom{(2)}{ij} = (z_i \quad i) - ! i;$$

The remaining cos-sin terms equal zero.

E PREDICTIVE DISTRIBUTION

To form predictive distributions, we assume that the latent posterior is approximated by $q(\mathbf{w}^{(r)}|\mathbf{X}; \cdot)$ as in Equation (20).

E.1 Predictive Intensity Distribution

For some $\mathbf{x} \in X$, the predictive distribution of $f(\mathbf{x})$ can be deduced from Equations (14) and (20) in the main text to be

$$f(\mathbf{x})|\mathbf{X}; \cdot \sim N(\hat{\mathbf{w}}^{(r)} \cdot^{(r)}(\mathbf{x}); \hat{\mathbf{w}}^{(r)} \cdot^{(r)}(\mathbf{x})) \quad (47)$$

where

$$\hat{\mathbf{w}}^{(r)} \cdot^{(r)}(\mathbf{x}) := \hat{\mathbf{w}}^{(r)\top} \cdot^{(r)}(\mathbf{x}) \quad (48)$$

and

$$\hat{\mathbf{w}}^{(r)} \cdot^{(r)}(\mathbf{x}) \cdot^{(r)}(\mathbf{x}) := \hat{\mathbf{w}}^{(r)\top}(\mathbf{x}) \mathbf{Q} \cdot^{(r)}(\mathbf{x}): \quad (49)$$

Given $\hat{\mathbf{w}}^{(r)} \cdot^{(r)}(\mathbf{x}) = (f(\mathbf{x}) + \cdot^2)$ and Equation (47) above, we can also derive the predictive distribution of the intensity function

$$\hat{\mathbf{w}}^{(r)} \cdot^{(r)}(\mathbf{x})|\mathbf{X}; \cdot \sim \text{Gamma}(a(\hat{\mathbf{w}}^{(r)} \cdot^{(r)}(\mathbf{x})); b(\hat{\mathbf{w}}^{(r)} \cdot^{(r)}(\mathbf{x}))) \quad (50)$$

with parameters

$$a(\hat{\mathbf{w}}^{(r)} \cdot^{(r)}(\mathbf{x})) = \frac{(\hat{\mathbf{w}}^{(r)} \cdot^{(r)}(\mathbf{x}))^2 + (\hat{\mathbf{w}}^{(r)} \cdot^{(r)}(\mathbf{x}))^2}{2(\hat{\mathbf{w}}^{(r)} \cdot^{(r)}(\mathbf{x}))^2(2(\hat{\mathbf{w}}^{(r)} \cdot^{(r)}(\mathbf{x}))^2 + (\hat{\mathbf{w}}^{(r)} \cdot^{(r)}(\mathbf{x}))^2)} \text{ and} \quad (51)$$

$$b(\hat{\mathbf{w}}^{(r)} \cdot^{(r)}(\mathbf{x})) = \frac{(\hat{\mathbf{w}}^{(r)} \cdot^{(r)}(\mathbf{x}))^2 + (\hat{\mathbf{w}}^{(r)} \cdot^{(r)}(\mathbf{x}))^2}{2(\hat{\mathbf{w}}^{(r)} \cdot^{(r)}(\mathbf{x}))^2(2(\hat{\mathbf{w}}^{(r)} \cdot^{(r)}(\mathbf{x}))^2 + (\hat{\mathbf{w}}^{(r)} \cdot^{(r)}(\mathbf{x}))^2)}. \quad (52)$$

E.2 Predictive Expected Log-likelihood

For a training set $\mathbf{X} = \{\mathbf{x}_i\}_{i=1}^N$ and an held-out test set $\mathbf{X} = \{\mathbf{x}_i\}_{i=1}^N$, we can derive an approximation for the expected predictive log-likelihood

$$E[\log p(\mathbf{X}|\mathbf{X})] \approx E_{\mathbf{w}^{(r)}} \int_{\mathbf{X}} (\mathbf{w}^{(r)\top} \cdot^{(r)}(\mathbf{x}) + \cdot^2) d\mathbf{x} + \sum_{i=1}^N E_{\mathbf{w}^{(r)}} \log(\mathbf{w}^{(r)\top} \cdot^{(r)}(\mathbf{x}_i) + \cdot^2)$$

where $\mathbf{w}^{(r)} \sim q(\mathbf{w}^{(r)}|\mathbf{X}; \cdot)$.

The integral term can be solved as

$$\begin{aligned} & E_{\mathbf{w}^{(r)}} \int_{\mathbf{X}} (\mathbf{w}^{(r)\top} \cdot^{(r)}(\mathbf{x}) + \cdot^2) d\mathbf{x} \\ &= \int_{\mathbf{X}} E_{\mathbf{w}^{(r)}} [\mathbf{w}^{(r)\top} \cdot^{(r)}(\mathbf{x}) + \cdot^2] d\mathbf{x} + \int_{\mathbf{X}} \text{Var}[\mathbf{w}^{(r)\top} \cdot^{(r)}(\mathbf{x})] d\mathbf{x} \\ &= \int_{\mathbf{X}} \hat{\mathbf{w}}^{(r)\top} \cdot^{(r)}(\mathbf{x}) \hat{\mathbf{w}}^{(r)} \cdot^{(r)}(\mathbf{x}) d\mathbf{x} + 2 \int_{\mathbf{X}} \hat{\mathbf{w}}^{(r)\top} \cdot^{(r)}(\mathbf{x}) d\mathbf{x} + \int_{\mathbf{X}} \cdot^2 d\mathbf{x} + \int_{\mathbf{X}} \hat{\mathbf{w}}^{(r)\top}(\mathbf{x}) \mathbf{Q} \cdot^{(r)}(\mathbf{x}) d\mathbf{x} \\ &= \text{tr}(\hat{\mathbf{W}}^{(r)} \hat{\mathbf{W}}^{(r)\top} + \mathbf{Q}) \int_{\mathbf{X}} \cdot^{(r)}(\mathbf{x}) \cdot^{(r)}(\mathbf{x}) d\mathbf{x} + 2 \int_{\mathbf{X}} \hat{\mathbf{w}}^{(r)\top} \cdot^{(r)}(\mathbf{x}) d\mathbf{x} + \int_{\mathbf{X}} \cdot^2 d\mathbf{x} \\ &= \hat{\mathbf{W}}^{(r)\top} \mathbf{M}^{(r)} \hat{\mathbf{W}}^{(r)} + \text{tr}(\mathbf{Q} \mathbf{M}^{(r)}) + 2 \int_{\mathbf{X}} \hat{\mathbf{w}}^{(r)\top} \cdot^{(r)}(\mathbf{x}) d\mathbf{x} + \int_{\mathbf{X}} \cdot^2 d\mathbf{x} \end{aligned}$$

where $\mathbf{M}^{(r)}$ and $\mathbf{m}^{(r)}$ are defined in Proposition 3.3. Note that we used the cyclical property of the trace in the last two lines. We also used the Tonelli's theorem in the first line to reverse the ordering of the integration over the positive integrand $(\mathbf{w}^{(r)\top} \cdot^{(r)}(\mathbf{x}) + \cdot^2)q(\mathbf{w}^{(r)})$.

The sum-of-expectations can also be expressed analytically. It takes of form

$$\sum_i E[\log z_i^2] \quad \text{where} \quad z_i \sim N(\cdot; \cdot) \quad (53)$$

with

$$z_i := \mathbf{w}^{(r)\top} \mathbf{x}_i + \epsilon_i \quad \text{and} \quad \epsilon_i := \mathbf{x}_i^\top \mathbf{Q} \mathbf{x}_i + \epsilon_i \quad (54)$$

Following Lloyd et al. (2015, section 4.3), each summand can be expressed as

$$\mathbb{E}[\log z_i^2] = G\left(\frac{i}{2}\right) + \log \frac{i}{2} + C$$

where $G(\cdot)$ is defined as

$$G(z) = 2z \prod_{j=0}^{\infty} \frac{j! z^j}{(2j)(1-2j)} \quad (55)$$

with $(\cdot)_j$ being the rising Pochhammer series. The constant $C = 0.57721566$ is the Euler Mascheroni constant. $G(\cdot)$ can in practice be evaluated using a large multi-resolution look-up table of pre-computed values. Accurate evaluation can be obtained by linear interpolation of the values from the table.

F MARGINAL LIKELIHOOD DERIVATIVES

In Equation (23) the marginal likelihood is expressed as

$$\begin{aligned} \log p(\mathbf{X}_j) &= \log p(\mathbf{w}^{(r)}; \mathbf{X}_j) - \log p(\mathbf{w}^{(r)}; \mathbf{X}; \epsilon) \\ &= \log p(\mathbf{w}^{(r)}; \mathbf{X}_j) - \log q(\mathbf{w}^{(r)}; \mathbf{X}; \epsilon) \\ &= \mathbf{w}^{(r)\top} \mathbf{M}^{(r)} \mathbf{w}^{(r)} + \sum_{i=1}^N \log(\mathbf{w}^{(r)\top} \mathbf{x}_i)^2 - \frac{1}{2} \mathbf{w}^{(r)\top} \mathbf{Q} \mathbf{w}^{(r)} + \frac{1}{2} \log j \mathbf{Q} + C \end{aligned}$$

for some constant C , where we assume $\epsilon = 0$ without loss of generality.

Marginal Likelihood Derivatives We now compute the gradient with respect to the hyperparameters θ_j . Using the chain rule,

$$\frac{\partial}{\partial \theta_j} \log p(\mathbf{X}_j) = \frac{\partial \log p(\mathbf{X}_j)}{\partial \mathbf{w}_j} \Big|_{\text{explicit}} + \sum_{j=1}^N \frac{\partial \log p(\mathbf{X}_j)}{\partial \mathbf{w}_j} \frac{\partial \mathbf{w}_j}{\partial \theta_j} \quad (56)$$

The first term of Equation (56) can be solved as

$$\frac{\partial \log p(\mathbf{X}_j)}{\partial \theta_j} \Big|_{\text{explicit}} = \mathbf{w}^{(r)\top} \frac{\partial \mathbf{M}^{(r)}}{\partial \theta_j} \mathbf{w}^{(r)} + 2 \sum_{i=1}^N \frac{\mathbf{w}^{(r)\top} \mathbf{x}_i}{\mathbf{w}^{(r)\top} \mathbf{x}_i} \frac{\partial \mathbf{x}_i}{\partial \theta_j} + \frac{1}{2} \text{tr}(\mathbf{Q}^{-1} \frac{\partial \mathbf{Q}}{\partial \theta_j}) \quad (57)$$

where $\mathbf{Q}^{-1} = \frac{\partial^2}{\partial \mathbf{w}^{(r)} \partial \mathbf{w}^{(r)}} \log p(\mathbf{w}^{(r)}; \mathbf{X}; \epsilon) \Big|_{\mathbf{w}^{(r)} = \mathbf{w}^{(r)}}$ is the precision matrix expressed in Section 4.1. The last term of Equation (57) can be expressed as

$$\begin{aligned} \frac{1}{2} \text{tr}(\mathbf{Q}^{-1} \frac{\partial \mathbf{Q}}{\partial \theta_j}) &= \frac{1}{2} \text{tr}(\frac{\partial \mathbf{Q}^{-1}}{\partial \theta_j} \mathbf{Q}) \\ &= \text{tr} \left(\frac{\partial \mathbf{M}^{(r)}}{\partial \theta_j} \mathbf{Q} + \frac{\partial}{\partial \theta_j} \sum_{i=1}^N \frac{\mathbf{x}_i \mathbf{x}_i^\top}{(\mathbf{w}^{(r)\top} \mathbf{x}_i)^2} \mathbf{Q} \right) \\ &= \text{tr} \left(\frac{\partial \mathbf{M}^{(r)}}{\partial \theta_j} \mathbf{Q} + \sum_{i=1}^N \frac{2}{(\mathbf{w}^{(r)\top} \mathbf{x}_i)^2} \mathbf{x}_i \mathbf{x}_i^\top \mathbf{Q} \frac{\partial \mathbf{x}_i}{\partial \theta_j} \right) \\ &\quad + 2 \sum_{i=1}^N \frac{\mathbf{w}^{(r)\top} \mathbf{x}_i}{(\mathbf{w}^{(r)\top} \mathbf{x}_i)^3} \mathbf{x}_i \mathbf{x}_i^\top \mathbf{Q} \mathbf{x}_i \end{aligned} \quad (58)$$

Algorithm 1 Compute the log marginal likelihood derivatives ^a.

- 1: **input:** \mathbf{X} (inputs), spectral locations and hyper-parameters $\tilde{\mathbf{w}}^{(r)}$, precision matrix \mathbf{Q}^{-1} , r "integral" matrix $\mathbf{M}^{(r)}$, N features matrix $\mathbf{r}^{(r)}(\mathbf{X})$
- 2: compute $\mathbf{r}^{(r)}$, the $\dim(\tilde{\mathbf{w}}) \times r$ tensor of partial derivatives of $\mathbf{M}^{(r)}$ with respect to $\tilde{\mathbf{w}}$
- 3: compute $\mathbf{r}^{(r)}(\mathbf{X})$, the $\dim(\tilde{\mathbf{w}}) \times N \times r$ tensor of partial derivatives of $\mathbf{r}^{(r)}(\mathbf{X})$ with respect to $\tilde{\mathbf{w}}$
- 4: $\hat{\mathbf{f}} := \mathbf{r}^{(r)}(\mathbf{X})\tilde{\mathbf{w}}^{(r)}$
- 5: $\mathbf{L} = \text{cholesky}(\mathbf{Q}^{-1})$. Solve $\mathbf{L}\mathbf{L}^\top = \mathbf{Q}^{-1}$
- 6: $\mathbf{C}_\cdot := \mathbf{L} \mathbf{r}^{(r)}(\mathbf{X})$
- 7: $\mathbf{C}_{r\cdot} := \mathbf{L} \mathbf{r}^{(r)}(\mathbf{X})$
- 8: $\mathbf{r}_\cdot := (\mathbf{C}_\cdot \ \mathbf{C}_{r\cdot}) \mathbf{1}_r$
- 9: $\mathbf{R}_{r\cdot} = (\mathbf{C}_\cdot \ \mathbf{C}_{r\cdot}) \mathbf{1}_r$
- 10: $\mathbf{S} := \mathbf{r}^{(r)}(\mathbf{X})\tilde{\mathbf{w}}^{(r)}$
- 11: $\mathbf{E} := (\mathbf{S} \ \text{diag}(\hat{\mathbf{f}}^2)) \mathbf{r}^{(r)}(\mathbf{X})$
- 12: $\mathbf{s}_{integral} = \tilde{\mathbf{w}}^{(r)\top} \mathbf{r}^{(r)} \mathbf{M}^{(r)} \tilde{\mathbf{w}}^{(r)}$
- 13: $\mathbf{s}_{data} = 2 (\mathbf{S} \ \text{diag}(\hat{\mathbf{f}}^2)) \mathbf{1}_N$
- 14: $\mathbf{s}_1 = \text{tr}(\mathbf{L}^\top \mathbf{r}^{(r)} (\mathbf{L} \mathbf{r}^{(r)} \mathbf{M}^{(r)}))$
- 15: $\mathbf{s}_2 = 2 (\mathbf{R}_{r\cdot} \ \text{diag}(\hat{\mathbf{f}}^2)) \mathbf{1}_N$
- 16: $\mathbf{s}_3 = 2 (\mathbf{S} \ \text{diag}(\hat{\mathbf{f}}^3)) \mathbf{r}_\cdot$
- 17: $\mathbf{v} := \text{diag}(\hat{\mathbf{f}}^2) \mathbf{r}^{(r)}(\mathbf{X}) \mathbf{1}_N^\top \mathbf{E}$. Equation (60)
- 18: $\frac{d\mathbf{p}}{d\tilde{\mathbf{w}}} := \mathbf{s}_{data} \ \mathbf{s}_{integral} \ \mathbf{s}_1 \ \mathbf{s}_2 + \mathbf{s}_3$. Equation (57)
- 19: $\frac{d\mathbf{p}}{d\tilde{\mathbf{w}}} := 2 \mathbf{r}^{(r)}(\mathbf{X})^\top \text{diag}(\hat{\mathbf{f}}^3) \mathbf{r}_\cdot$. Equation (59)
- 20: $\frac{d\mathbf{W}}{d\tilde{\mathbf{w}}} := 2 (\mathbf{L} \mathbf{r}^{(r)} (\mathbf{L}^\top \mathbf{r}^{(r)} \mathbf{M}^{(r)} \tilde{\mathbf{w}}^{(r)} \ \mathbf{v}))$. Equation (61)
- 21: $\mathbf{g} = \frac{d\mathbf{p}}{d\tilde{\mathbf{w}}} + \frac{d\mathbf{W}}{d\tilde{\mathbf{w}}} \frac{d\mathbf{p}}{d\tilde{\mathbf{w}}}$. Equation (56) **return** \mathbf{g} ($\dim(\tilde{\mathbf{w}})$ -vector of partial derivatives)

^aWe assume basic operations on $3d$ tensors are performed over its last two dimensions and repeated over its first one. $*$ denotes the element-wise multiplication. $A \setminus B$ where A is a triangular $r \times r$ matrix and B is a $\dots \times r$ matrix or tensor is performed over the last dimension of B (i.e. r) and repeated over its firsts dimensions.

The $\frac{\partial \log p(\mathbf{X}_j)}{\partial \tilde{\mathbf{w}}_j}$ terms of Equation (56) is

$$\begin{aligned}
 \frac{\partial \log p(\mathbf{X}_j)}{\partial \tilde{\mathbf{w}}_j} &= \frac{\partial \log p(\tilde{\mathbf{w}}^{(r)}; \mathbf{X}_j)}{\partial \tilde{\mathbf{w}}_j} + \frac{1}{2} \text{tr}(\mathbf{Q}^{-1} \frac{\partial \mathbf{Q}}{\partial \tilde{\mathbf{w}}_j}) \\
 &= \frac{1}{2} \text{tr}(\frac{\partial \mathbf{Q}^{-1}}{\partial \tilde{\mathbf{w}}_j} \mathbf{Q}) \\
 &= \text{tr} \frac{\partial}{\partial \tilde{\mathbf{w}}_j} \sum_{i=1}^N \frac{\mathbf{r}^{(r)}(\mathbf{x}_i) \mathbf{r}^{(r)\top}(\mathbf{x}_i)}{(\tilde{\mathbf{w}}^{(r)\top} \mathbf{r}^{(r)}(\mathbf{x}_i))^2} \mathbf{Q} \\
 &= 2 \sum_{i=1}^N \frac{\mathbf{r}^{(r)}(\mathbf{x}_i)}{(\tilde{\mathbf{w}}^{(r)\top} \mathbf{r}^{(r)}(\mathbf{x}_i))^3} \mathbf{r}^{(r)}(\mathbf{x}_i) \mathbf{Q} \mathbf{r}^{(r)\top}(\mathbf{x}_i)
 \end{aligned} \tag{59}$$

where, in the first line, we have imposed stationarity using $\mathbf{r}^{(r)} \log p(\tilde{\mathbf{w}}^{(r)}; \mathbf{X}_j) |_{\tilde{\mathbf{w}}^{(r)} = \tilde{\mathbf{w}}^{(r)}} = 0$.

We finally express the last terms $\frac{\partial \tilde{\mathbf{w}}^{(r)}}{\partial \tilde{\mathbf{w}}_i}$ of Equation (56). From the expression of $\tilde{\mathbf{w}}^{(r)}$ in Equation (21) in the main text. We

obtain

$$\begin{aligned}
 \frac{\partial \hat{\mathbf{w}}^{(r)}}{\partial \mathbf{w}_j} &= \frac{\partial}{\partial \mathbf{w}_j} \left(\mathbf{M}^{(r)} + \frac{1}{2} I_r \sum_{i=1}^N \frac{\mathbf{x}_i^{(r)} \mathbf{x}_i^{(r)^\top}{\mathbf{w}^{(r)}}}{\|\mathbf{w}^{(r)}\|^2} \right) \\
 &= \frac{\partial}{\partial \mathbf{w}_j} \left(\mathbf{M}^{(r)} + \frac{1}{2} I_r \sum_{i=1}^N \frac{\mathbf{x}_i^{(r)} \mathbf{x}_i^{(r)^\top}{\mathbf{w}^{(r)}}}{\|\mathbf{w}^{(r)}\|^2} \right) + \frac{\partial}{\partial \hat{\mathbf{w}}^{(r)}} \sum_{i=1}^N \frac{\mathbf{x}_i^{(r)} \mathbf{x}_i^{(r)^\top}{\mathbf{w}^{(r)}}}{\|\mathbf{w}^{(r)}\|^2} \frac{\partial \hat{\mathbf{w}}^{(r)}}{\partial \mathbf{w}_j} \\
 &= \left(\mathbf{M}^{(r)} + \frac{1}{2} I_r \sum_{i=1}^N \frac{\mathbf{x}_i^{(r)} \mathbf{x}_i^{(r)^\top}{\mathbf{w}^{(r)}}}{\|\mathbf{w}^{(r)}\|^2} \right) \frac{\partial \hat{\mathbf{w}}^{(r)}}{\partial \mathbf{w}_j} + \sum_{i=1}^N \frac{\mathbf{x}_i^{(r)} \mathbf{x}_i^{(r)^\top}{\mathbf{w}^{(r)}}}{\|\mathbf{w}^{(r)}\|^2} \frac{\partial \hat{\mathbf{w}}^{(r)}}{\partial \mathbf{w}_j}
 \end{aligned}$$

where

$$\hat{\mathbf{v}} := \frac{\partial}{\partial \mathbf{w}_j} \sum_{i=1}^N \frac{\mathbf{x}_i^{(r)} \mathbf{x}_i^{(r)^\top}{\mathbf{w}^{(r)}}}{\|\mathbf{w}^{(r)}\|^2} = \sum_{i=1}^N \frac{\frac{\partial}{\partial \mathbf{w}_j} (\mathbf{x}_i^{(r)} \mathbf{x}_i^{(r)^\top)} \mathbf{w}^{(r)}}{\|\mathbf{w}^{(r)}\|^2} - \sum_{i=1}^N \frac{\mathbf{x}_i^{(r)} \mathbf{x}_i^{(r)^\top}{\mathbf{w}^{(r)}}}{\|\mathbf{w}^{(r)}\|^4} \mathbf{w}^{(r)} \quad (60)$$

Thus

$$\begin{aligned}
 \frac{\partial \hat{\mathbf{w}}^{(r)}}{\partial \mathbf{w}_j} &= \left(\mathbf{M}^{(r)} + \frac{1}{2} I_r + \sum_{i=1}^N \frac{\mathbf{x}_i^{(r)} \mathbf{x}_i^{(r)^\top}{\mathbf{w}^{(r)}}}{\|\mathbf{w}^{(r)}\|^2} \right) \frac{\partial \hat{\mathbf{w}}^{(r)}}{\partial \mathbf{w}_j} + \hat{\mathbf{v}} \\
 &= 2\mathbf{Q} \frac{\partial \hat{\mathbf{w}}^{(r)}}{\partial \mathbf{w}_j} + \hat{\mathbf{v}}
 \end{aligned} \quad (61)$$

Implementation Details The implementation of the marginal log likelihood partial derivatives with respect to the hyperparameters $\tilde{\mathbf{w}}$ is shown in Algorithm 1.

G ALGORITHMS

We present in Algorithm 2 the standard process to compute the mode $\hat{\mathbf{w}}^{(r)}$ and the hyperparameters $\tilde{\mathbf{w}}$.

Algorithm 2 Compute the mode $\hat{\mathbf{w}}^{(r)}$ and the hyperparameters

- 1: **input:** data \mathbf{X}
 - 2: initialize \mathbf{w}_0 and $\hat{\mathbf{w}}_0^{(r)}$
 - 3: **for** $t = 1; \dots; T$ **do**
 - 4: compute $\mathbf{M}^{(r)}$ and $\mathbf{x}^{(r)}(\mathbf{X})$
 - 5: $\hat{\mathbf{w}}_t^{(r)} := \text{mode}(p(\mathbf{w}^{(r)} | \mathbf{X}; \tilde{\mathbf{w}}); \hat{\mathbf{w}}_{t-1}^{(r)})$. locate posterior mode using Equation (22) with initial value $\hat{\mathbf{w}}_{t-1}^{(r)}$
 - 6: $\hat{\mathbf{f}} := \mathbf{x}^{(r)}(\mathbf{X}) \hat{\mathbf{w}}_t^{(r)}$
 - 7: $\mathbf{V} := \text{diag}(\hat{\mathbf{f}}^{-1}) \mathbf{x}^{(r)}(\mathbf{X})$
 - 8: $\mathbf{Q}^{-1} := 2\mathbf{M}^{(r)} + I_r + 2\mathbf{V}^\top \mathbf{V}$. precision matrix $\mathbf{Q}^{-1} = r_{\mathbf{w}^{(r)}}^2 \log p(\mathbf{w}^{(r)} | \mathbf{X}; \tilde{\mathbf{w}})_{\mathbf{w}^{(r)} = \hat{\mathbf{w}}_t^{(r)}}$
 - 9: compute the gradient \mathbf{g} from Algorithm 1
 - 10: $\tilde{\mathbf{w}}_t$ update($\tilde{\mathbf{w}}_{t-1}; \mathbf{g}$)
 - 11: **end for**
 - 12: **return** $\hat{\mathbf{w}}_T^{(r)}$ (mode) and $\tilde{\mathbf{w}}_T$ (hyperparameters)
-

As an alternative (Section 4.2), we assume $\frac{\partial \hat{\mathbf{w}}^{(r)}}{\partial \mathbf{w}_j} = 0$ and alternate independent update for the mode in Equation (22) and the hyperparameters in Equation (24), which yields Algorithm 3.

Algorithm 3 Compute the mode $\hat{\mathbf{w}}^{(r)}$ and the hyperparameters with independent updates

```

1: input: data  $\mathbf{X}$ 
2: initialize  $\mathbf{w}_0$  and  $\hat{\mathbf{w}}_0^{(r)}$ 
3: for  $t = 1; \dots; T$  do
4:   compute  $\mathbf{M}^{(r)}$  and  $\mathbf{V}^{(r)}(\mathbf{X})$ 
5:    $\hat{\mathbf{w}}_t^{(r)} = \hat{\mathbf{w}}_t^{(r)}$  . update the posterior mode using one iteration of Equation (22)
6:    $\hat{\mathbf{f}} := \mathbf{V}^{(r)}(\mathbf{X})\hat{\mathbf{w}}_t^{(r)} + \mathbf{y}$ 
7:    $\mathbf{V} := \text{diag}(\hat{\mathbf{f}}^{-1})\mathbf{V}^{(r)}(\mathbf{X})$ 
8:    $\mathbf{Q}^{-1} := 2\mathbf{M}^{(r)} + I_r + 2\mathbf{V}^{-1}\mathbf{V}$  . precision matrix
9:   compute the gradient  $\frac{d\mathcal{L}}{d\mathbf{p}}$  from Algorithm 1 (line 18)
10:   $\mathbf{p}_t$  update( $\mathbf{p}_{t-1}; \frac{d\mathcal{L}}{d\mathbf{p}}$ )
11: end for
12: return  $\hat{\mathbf{w}}_T^{(r)}$  (mode) and  $\mathbf{p}_T$  (hyperparameters)
    
```

H SUPPLEMENTARY FIGURES

Figure 5 (below) shows the average test expected log-likelihood as a function of the number of basis functions or inducing points for *bei* data set and *Taxi* data set across models.

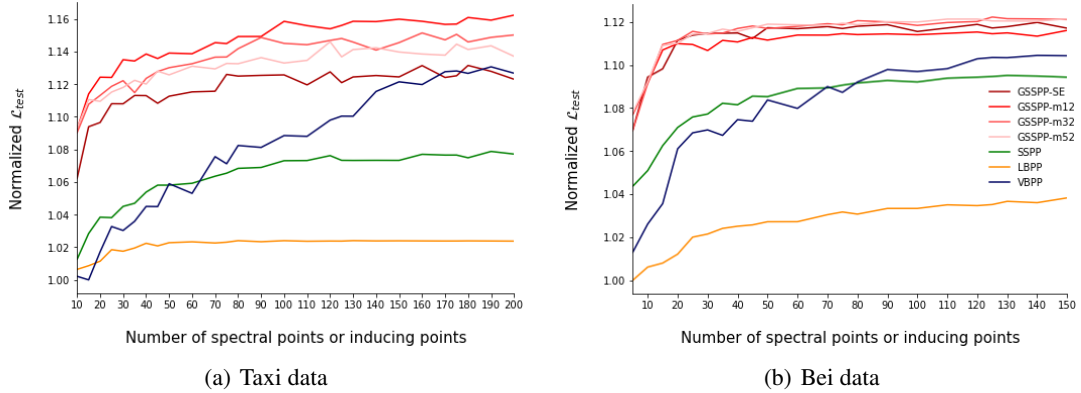


Figure 5: Average normalized test expected log-likelihood (L_{test}) of the different methods on 2D real data, as a function of the number of spectral points or inducing points.

