

Learning Compositional Sparse Gaussian Processes with a Shrinkage Prior

Anh Tong¹, Toan Tran², Hung Bui², Jaesik Choi^{3,4}

¹ Ulsan National Institute of Science and Technology

² VinAI Research

³ Korea Advanced Institute of Science and Technology

⁴ INEEJI

anhth@unist.ac.kr, {v.toantm3,v.hungbh1}@vinai.io, jaesik.choi@kaist.ac.kr

Abstract

Choosing a proper set of kernel functions is an important problem in learning Gaussian Process (GP) models since each kernel structure has different model complexity and data fitness. Recently, automatic kernel composition methods provide not only accurate prediction but also attractive interpretability through search-based methods. However, existing methods suffer from slow kernel composition learning. To tackle large-scaled data, we propose a new sparse approximate posterior for GPs, MultiSVGP, constructed from groups of inducing points associated with individual additive kernels in compositional kernels. We demonstrate that this approximation provides a better fit to learn compositional kernels given empirical observations. We also provide theoretical justification on error bound when compared to the traditional sparse GP. In contrast to the search-based approach, we present a novel probabilistic algorithm to learn a kernel composition by handling the sparsity in the kernel selection with Horseshoe prior. We demonstrate that our model can capture characteristics of time series with significant reductions in computational time and have competitive regression performance on real-world data sets.

1 Introduction

Recently, there are numerous advancements in automating model learning with statistical methods. There are still many challenging problems including making the automatic procedure efficient and applying methods on large-scaled data.

The Automatic Statistician framework (Grosse et al. 2012; Duvenaud et al. 2013; Lloyd et al. 2014; Ghahramani 2015; Steinruecken et al. 2019) aims to address challenges on automating model discovery. The framework adopts search procedures over a space of models, enumerating possible compositional Gaussian Process (GP) kernel functions generated from base ones and compositional grammar rules. Learned models can produce human-readable explanations as dissecting small explainable components in the compositional kernel models. However, existing methods suffer from slow kernel composition learning due to the huge search space.

Recently, a differential compositional learning inspired by deep neural networks (Sun et al. 2018) is proposed. Although this model obtains expressive kernel functions for GPs and consequently achieves good predictive performances, it is hard to interpret compared to existing Bayesian kernel learning methods, e.g. (Lloyd et al. 2014).

This paper presents a new kernel composition learning method which encourages to learn a sparse composition of kernels with a shrinkage prior, Horseshoe prior, which is proven to be effective in learning sparse signal (Carvalho, Polson, and Scott 2009; Bhadra et al. 2019). To retain the model interpretability, we formulate compositional kernels as the sum of individual additive kernels which can be explained in natural language.

To scale up the model for large-scaled data, we introduce a new approximate posterior for GP, Multi-inducing Sparse Variational GP (MultiSVGP). Existing sparse inducing GP (Snelson and Ghahramani 2006; Hensman, Fusi, and Lawrence 2013) methods give an approximation for GP for a kernel function in general with a single set of inducing points. However, we can further improve this approach specifically for additive compositional kernels by considering a group of inducing points and assigning an individual member of inducing points responsible for an additive kernel in our approximating posterior. We justify that our approximation with compositional kernels produces a better error bound than the sparse inducing GP. In experiments, we demonstrate that our models can capture appropriate kernel structures for time series from small-scaled data sets to large-scaled data sets. In general, our model runs faster than existing compositional kernel learning methods (Lloyd et al. 2014; Kim and Teh 2018) five to twenty times while maintaining similar accuracy. We also show that our model to have competitive extrapolation performance in regression tasks as well as improve additive GPs (Duvenaud, Nickisch, and Rasmussen 2011) with our kernel selection.

The paper offers the following contributions: (1) a new accurate GP approximation for additive compositional kernels; (2) kernel selection using Horseshoe prior so that the learned models can capture the inductive kernel structure in data and remain interpretable.

2 Related work

There is a large body of work (Duvenaud et al. 2013; Lloyd et al. 2014; Kim and Teh 2018) establishing the foundation of model discovery for Gaussian processes. (Grosse et al. 2012) presents work on unsupervised learning for the case of matrix decomposition. Then, (Duvenaud et al. 2013; Lloyd et al. 2014; Ghahramani 2015; Hwang, Tong, and Choi 2016; Tong and Choi 2019) extend to supervised settings with Gaussian process (GP) models. (Sun et al. 2018) build complex kernels under network architectures having an additive layer where any two kernels are summed and followed by a product layer where kernels are multiplied. (Kim and Teh 2018) adopt search procedure but smartly avoid the learning model by finding bounds of likelihood. However, the search is still time-consuming and is done heuristically greedy manner. A complete review as well as a guideline for automatic systems can be found in (Steinruecken et al. 2019). While inheriting the spirits of existing work on kernel compositions, this paper focuses on scaling up the system in terms of data size and efficient model selection.

A recent work (Teng et al. 2020) shares some similarity to our work. The paper presents a probabilistic approach to select among models with a Softmax-like assumption for choosing a model. The main idea in this paper is to select a single model out of a manual fixed set of candidate models. Our model generating kernels combinatorially considers a bigger model space than that of (Teng et al. 2020). On the other hand, models (Teng et al. 2020) use single inducing points for compositional kernels. Its sparse GP approximation can be limited compared to our MultiSVGP. Another work (Malkomes, Schaff, and Garnett 2016) attempts to extract a model out of candidate ones using surrogate models, e.g., Bayesian optimization. The surrogate models are based on the distance between GP models. A recent approach on learning implicit kernel (Tompkins et al. 2019) seeks for the representation in their corresponding spectral domain. Exploiting the spectral representation is previously used in (Wilson and Adams 2013). However, the question of how to make these methods interpretable left unanswered.

There is existing work proposing probabilistic priors, e.g., spike-and-slab prior, on multi-task GP (Titsias and Lázaro-Gredilla 2011). However, the approach does not scale well with the number of data and suffers from computational burden by the choice of probabilistic priors. Other work that is related but distinct includes (Senanayake, Tompkins, and Ramos 2018; Paciorek and Schervish 2003; Tong and Choi 2016; Saad et al. 2019).

3 Background

This section presents reviews on building blocks of this paper including Gaussian process, the kernel selection framework, and shrinkage prior.

3.1 Gaussian Process

Gaussian Process (GP) is known as a nonparametric prior over function values (Rasmussen and Williams 2005). Formally, a GP is defined as a multivariate Gaussian distribution

over function value $f(\mathbf{x})$:

$$f(\mathbf{x}) \sim \mathcal{GP}(m(\mathbf{x}), k(\mathbf{x}, \mathbf{x}')),$$

where $m(\mathbf{x})$ is the mean function (usually is set as zero mean), and $k(\mathbf{x}, \mathbf{x}')$ is the kernel (covariance) function. Given data $\mathcal{D} = (\mathbf{X}, \mathbf{y})$, the predictive posterior distribution at \mathbf{x}_* is in a closed form:

$$\begin{aligned} f(\mathbf{x}_*) | \mathbf{X}, \mathbf{y} &\sim \mathcal{N}(\mu(\mathbf{x}_*), \sigma^2(\mathbf{x}_*)), \\ \mu(\mathbf{x}_*) &= k(\mathbf{x}_*, \mathbf{X})k(\mathbf{X}, \mathbf{X})^{-1}\mathbf{y}, \\ \sigma^2(\mathbf{x}_*) &= k(\mathbf{x}_*, \mathbf{x}_*) - k(\mathbf{x}_*, \mathbf{X})k(\mathbf{X}, \mathbf{X})^{-1}k(\mathbf{x}_*, \mathbf{X})^\top, \end{aligned}$$

where $k(\mathbf{x}_*, \mathbf{X})$ computes the covariance between function evaluations at \mathbf{x}_* and \mathbf{X} . This property of GP provides an efficient way to calibrate model uncertainties, becoming useful in several calibration methods, i.e. Bayesian optimization and Bayesian quadrature.

3.2 The Automatic Statistician

The Automatic Statistician or Automatic Bayesian Covariance Discovery (ABCD) aims to mimic and automate the process of statistical modeling (Grosse et al. 2012; Duvenaud et al. 2013; Lloyd et al. 2014; Ghahramani 2015). There are three main components in this framework: language of models, search procedure and report generations.

Language of models GP models are constructed by a grammar over kernels with a set of base kernels and kernel operators. The base kernels are: SE (squared exponential), LIN (linear), PER (periodic) (see Appendix A for details). The operators consist of $+$ (addition), \times (multiplication). As composed kernels get more complex, the corresponding generated models become more expressive to fit complex data.

Search procedure The search procedure is done in a greedy manner. That is, the language of models generates candidate models. Then, all of them are optimized by maximizing log likelihoods. A model is selected based on the trade-off between model and data complexity. This selected model is the input of the language of models to create new candidates.

Automatic generated explanation of models The compositional kernels resulted from the search procedure are transformed into a sum of products of base kernels. Each structural product of kernels is interpreted under natural-language expressions. All descriptions are gathered to produce a complete report with visualized plots and human-friendly analyses.

3.3 Sparse Variational Gaussian Process

The history of sparse Gaussian process methods dated back from the work of (Snelson and Ghahramani 2006). It can be considered as a sibling of Nyström approximation (Williams and Seeger 2001). The main idea of sparse Gaussian process is to introduce *pseudo inducing points*, \mathbf{u} , which are distributed jointly with Gaussian process latent variable \mathbf{f} under a Gaussian distribution. The number of inducing points, M , is much smaller than the number of data points, N . The computational complexity of sparse Gaussian process

is $\mathcal{O}(NM^2)$ for each learning iteration. There is a line of research on improving and understanding sparse Gaussian processes (Bauer, van der Wilk, and Rasmussen 2016; Bui, Yan, and Turner 2017; Burt, Rasmussen, and Van Der Wilk 2019).

Given a data set $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$, a sparse variational Gaussian process (SVGP) is defined by

$$\begin{aligned} f(\cdot) &\sim \mathcal{GP}(0, k(\cdot, \cdot)), \\ y_i | f, \mathbf{x}_i &\sim p(y_i | f(\mathbf{x}_i)), \end{aligned}$$

where $p(y_i | f(\mathbf{x}_i))$ is a likelihood, e.g., Gaussian, categorical. To approximate the posterior, the variational approach considers M inducing points $\mathbf{u} = \{u_i\}_{i=1}^M$ at locations $\{\mathbf{z}_i\}_{i=1}^M$, forming the variational distribution as

$$\begin{aligned} \mathbf{u} &\sim \mathcal{N}(\mathbf{m}, \mathbf{S}), \\ f(\cdot) | \mathbf{u} &\sim \mathcal{GP}(\mu(\cdot), \Sigma(\cdot, \cdot)), \end{aligned}$$

where the mean and covariance are obtained as

$$\begin{aligned} \mu(\cdot) &= \mathbf{k}_u(\cdot)^\top \mathbf{K}_{uu}^{-1} \mathbf{u}, \\ \Sigma(\cdot, \cdot) &= k(\cdot, \cdot) - \mathbf{k}_u(\cdot)^\top \mathbf{K}_{uu}^{-1} \mathbf{k}_u(\cdot), \end{aligned} \quad (1)$$

with $\mathbf{k}_u(\cdot) = [k(\mathbf{z}_i, \cdot)]_{i=1}^M$ and \mathbf{K}_{uu} is the covariance of \mathbf{u} .

The variational inference maximizes the evidence lower bound (ELBO) given as following (Hensman, Fusi, and Lawrence 2013)

$$\mathcal{L} = \sum_i \mathbb{E}_{q(f(\cdot))} [\log p(y_i | f(\mathbf{x}_i))] - \text{KL}[q(\mathbf{u}) || p(\mathbf{u})].$$

Here $q(f(\cdot))$ is a Gaussian, having mean as $\mathbf{k}_u(\cdot)^\top \mathbf{K}_{uu}^{-1} \mathbf{m}$ and covariance as $k(\cdot, \cdot) - \mathbf{k}_u(\cdot)^\top \mathbf{K}_{uu}^{-1} (\mathbf{S} - \mathbf{K}_{uu}) \mathbf{K}_{uu}^{-1} \mathbf{k}_u(\cdot)$. The objective \mathcal{L} can be optimized using stochastic gradient descent.

3.4 Shrinkage prior

In many model learning scenarios, we usually face the problem of sparse modeling for variable selection. Some prominent methods are proposed to tackle the problem, including Lasso regularization (Tibshirani 1996), spike and slab prior (Mitchell and Beauchamp 1988) and Horseshoe prior (Carvalho, Polson, and Scott 2009).

Spike and slab prior The spike and slab prior over \mathbf{w} is defined as

$$\begin{aligned} v_i &\sim \mathcal{N}(0, \sigma_w^2), \\ s_i &\sim \text{Bernoulli}(\pi_s), \\ w_i &= v_i s_i. \end{aligned}$$

This belongs to the two-group models. That is, the event $\{w_i = 0\}$ happens with probability $(1 - \pi_s)$; and nonzero w_i distributed according to a Gaussian prior with probability π_s . There is existing work using the prior for kernel learning (Titsias and Lázaro-Gredilla 2011).

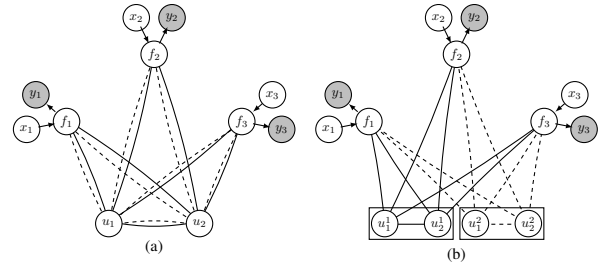


Figure 1: The graphical model of two models. Solid and dashed lines indicate the connections modeled by two different kernel function k_1 and k_2 . (a): Sparse inducing GP. The inducing points u_i is introduced as a proxy for the connections between f_i . (b): Our approach. Inducing points are grouped. Each group represents an individual kernel k_1 or k_2 .

Horseshoe prior The horseshoe prior (Carvalho, Polson, and Scott 2009) introduces a way to sample a sparse vector β as

$$\begin{aligned} \beta_i &\sim \mathcal{N}(0, \tau^2 \lambda_i^2), \quad i = 1 \dots m \\ \lambda_i &\sim \mathcal{C}^+(B), \quad i = 1 \dots m \\ \tau &\sim \mathcal{C}^+(A), \end{aligned}$$

where $\mathcal{C}^+(\cdot)$ is the half-Cauchy distribution, A and B are the scale parameters. Here, τ is the global shrinkage parameter, λ_i is the local shrinkage parameter. In contrast to the spike and slab prior, horseshoe prior is a continuous shrinkage one. It has Cauchy-like tails which allow signals to at large values. On the other hand, the infinite spike near zero keeps w_i around the origin. (Ghosh, Yao, and Doshi-Velez 2019) uses Horseshoe prior for the weight selection in Bayesian deep neural networks.

Compared to Horseshoe prior, the spike and slab prior exhibits a substantial computational burden as the dimension of sparse vectors increases.

4 Kernel selection with shrinkage prior

This section presents our main contributions: (1) kernel selection with Horseshoe prior and (2) our approximate GP for compositional kernels.

4.1 Kernel selection with Horseshoe prior

Consider the full GP model with kernel construction based on generative procedures:

$$f(\cdot) | \mathbf{w} \sim \mathcal{GP}(0, \tilde{k}(\cdot, \cdot)), \quad (2)$$

where $\tilde{k}(\mathbf{x}, \mathbf{x}') = \sum_{i=1}^m w_i^2 k_i(\mathbf{x}, \mathbf{x}')$ is constructed from m kernel functions $k_i(\mathbf{x}, \mathbf{x}')$. We introduce a probabilistic prior over the weights $\mathbf{w} = [w_{1:m}]$, $p(\mathbf{w})$ motivated from the Horseshoe prior. That is, we add the covariance term $k_i(\mathbf{x}, \mathbf{x}')$ to the step sampling β_i in Horseshoe generative procedure. This makes β_i equivalent to $f_i(\mathbf{x})$ in GP, i.e.,

$$\beta_i \sim \mathcal{N}(0, \tau^2 \lambda_i^2) \Rightarrow f_i(\mathbf{x}) \sim \mathcal{GP}(0, \tau^2 \lambda_i^2 k_i(\mathbf{x}, \mathbf{x}')).$$

When considering the multivariate normal distribution $\mathbf{f}_i \sim \mathcal{N}(\mathbf{0}, \tau^2 \lambda_i^2 \mathbf{K}_i)$ with kernel matrix \mathbf{K}_i computed from $k_i(\mathbf{x}, \mathbf{x}')$, the *multivariate* version of Horseshoe variable, $\beta_i \sim \mathcal{N}(\mathbf{0}, \tau^2 \lambda_i^2 \mathbf{I})$, is a special case of \mathbf{f}_i when \mathbf{K}_i is the identity matrix. This generalization is natural, equipping the sparsity among $\{f_i(\mathbf{x})\}_{i=1}^m$. Denoting $w_i^2 = \tau^2 \lambda_i^2$ and assuming that $\{f_i(\mathbf{x})\}_{i=1}^m$ are mutually independent, we can get $f(\mathbf{x}) = \sum_{i=1}^m f_i(\mathbf{x}) \sim \mathcal{GP}(0, \tilde{k}(\mathbf{x}, \mathbf{x}'))$.

The assumption on sparsity among kernel functions k_i encourages simple kernels which agree with model selection principles like Occam's razor in (Rasmussen and Ghahramani 2001) and BIC in (Lloyd et al. 2014).

4.2 Multi-inducing sparse Gaussian process

To motivate the proposed approach, we will first provide a naive model directly obtained from the sparse Gaussian process. Then, we present our main model.

SVGP with compositional kernels Given inducing points \mathbf{u} , we formulate the corresponding sparse GP model as

$$\begin{aligned} f(\cdot) | \mathbf{w}, \mathbf{u} &\sim \mathcal{GP}(\tilde{\mu}(\cdot), \tilde{\Sigma}(\cdot, \cdot)), \\ \tilde{\mu}(\cdot) &= \tilde{\mathbf{k}}_{\mathbf{u}}^\top(\cdot) \tilde{\mathbf{K}}_{\mathbf{u}\mathbf{u}}^{-1} \mathbf{u}, \\ \tilde{\Sigma}(\cdot, \cdot) &= \tilde{k}(\cdot, \cdot) - \tilde{\mathbf{k}}_{\mathbf{u}}^\top(\cdot) \tilde{\mathbf{K}}_{\mathbf{u}\mathbf{u}}^{-1} \tilde{\mathbf{k}}_{\mathbf{u}}(\cdot). \end{aligned} \quad (3)$$

Here, we denote that $\tilde{\mathbf{K}}_{\mathbf{u}\mathbf{u}} = \sum_{i=1}^m w_i^2 \mathbf{K}_{i\mathbf{u}\mathbf{u}}$, and $\mathbf{K}_{i\mathbf{u}\mathbf{u}}$ is the covariance of \mathbf{u} computed from kernel function $k_i(\cdot, \cdot)$. **Multi-inducing sparse variational GP** Given \mathbf{w} , we define the model via the combination of conditional posterior distributions:

$$f(\cdot) | \mathbf{U}, \mathbf{w} \sim \mathcal{GP} \left(\sum_{i=1}^m w_i \mu_i(\cdot; \mathbf{u}_i), \sum_{i=1}^m w_i^2 \Sigma_i(\cdot, \cdot; \mathbf{u}_i) \right), \quad (4)$$

where

$$\begin{aligned} \mu_i(\cdot) &= \mathbf{k}_{\mathbf{u}_i}(\cdot)^\top \mathbf{K}_{\mathbf{u}_i \mathbf{u}_i}^{-1} \mathbf{u}_i, \\ \Sigma_i(\cdot, \cdot) &= k_i(\cdot, \cdot) - \mathbf{k}_{\mathbf{u}_i}(\cdot)^\top \mathbf{K}_{\mathbf{u}_i \mathbf{u}_i}^{-1} \mathbf{k}_{\mathbf{u}_i}(\cdot). \end{aligned} \quad (5)$$

Here $k_{\mathbf{u}_i}(\cdot) = [k_i(\mathbf{u}_i, \cdot)]^\top$, and $\mathbf{K}_{\mathbf{u}_i \mathbf{u}_i}$ is the covariance of \mathbf{u}_i w.r.t. kernel k_i . For convenience, we omit the notation \mathbf{u}_i in μ_i and Σ_i .

In contrast to the model in Equation 3, we use m inducing groups of inducing points, $\mathbf{U} = \{\mathbf{u}_i\}_{i=1}^m$. Each \mathbf{u}_i is associated with a kernel function k_i , consisting of M_i inducing points $\mathbf{u}_i = [u_{1:M_i}^{(i)}]$ at inducing location $\mathbf{Z}_i = [\mathbf{z}_{1:M_i}^{(i)}] \in \mathcal{Z}_i$. The number of inducing points M_i is chosen to be much smaller than the size of data set N , i.e. ($M_i \ll N$).

Figure 1 is an illustrative comparison between SVGP and our proposed approach. In plain words, each member in inducing groups is responsible for a single kernel structural representation. We call this model as Multi-inducing Sparse Variational Gaussian Process (MultiSVGP).

Discussion It is obvious that the conditional distribution in Equation 4 is not equivalent to the conditional distribution of SVGP in Equation 3 since the inverse of the matrix sum is not equal to the sum of inverse matrices. Our proposed conditional distribution treats each kernel independently with

separate inducing points while the condition in SVGP contains correlations between the kernels which are often complicated under the matrix inverse operator.

Better fit We postulate that the combination of individual conditional Gaussian is still able to learn from data well comparing to other GP models. We justify this with a small experiment in which data are generated from a true model with kernel function $\text{SE}_1 + \text{SE}_2 + \text{PER}_1$. We then fit the data with full GP, SVGP model and our proposed approach. Figure 2 shows the posteriors corresponding to these models along with their Wasserstein-2 distance to the true model. We can observe the high quality of the posterior obtained from our assumption.

Connection to inter-domain variational Gaussian Process Inter-domain Gaussian process (Lázaro-Gredilla and Figueiras-Vidal 2009) helps finding the representative features which lie in difference domains. This work adopts the similar methodology to explain the proposed model.

Consider m Gaussian processes which are governed by different kernels:

$$g_i(\mathbf{x}) \sim \mathcal{GP}(0, k_i(\mathbf{x}, \mathbf{x}')), \quad i = 1 \dots m. \quad (6)$$

When applying a linear transformation over a GP, we can obtain a new GP (Lázaro-Gredilla and Figueiras-Vidal 2009). We consider the transformation

$$u_i(\cdot) = \int \phi_i(\mathbf{x}, \cdot) g_i(\mathbf{x}) d\mathbf{x}.$$

The choice of $\phi_i(\mathbf{x}, \mathbf{z})$ is a Dirac function with the information of which inducing point group $\tilde{\mathbf{z}}$ is in:

$$\phi_i(\mathbf{x}, \mathbf{z}) = \mathbb{I}\{\mathbf{z} \in \mathcal{Z}_i\} \delta(\mathbf{z} - \mathbf{x}).$$

Choosing Dirac delta function $\delta(\cdot)$ is similar to traditional sparse Gaussian process. Whereas $\mathbb{I}\{\mathbf{z} \in \mathcal{Z}_i\}$ provides the membership information of inducing points in the group. If $\mathbf{z} \in \mathcal{Z}_i$, then $\mathbb{I}\{\mathbf{z} \in \mathcal{Z}_i\} = 1$. We combine all $g_i(\mathbf{x})$ to get the model in Equation 2:

$$f(\mathbf{x}) = \sum_{i=1}^m w_i g_i(\mathbf{x}) \sim \mathcal{GP} \left(0, \sum_{i=1}^m w_i^2 k_i(\mathbf{x}, \mathbf{x}') \right).$$

Followed by (Lázaro-Gredilla and Figueiras-Vidal 2009), the corresponding (cross-) covariance between \mathbf{U} and \mathbf{f} can be obtained as

$$\begin{aligned} k_{\mathbf{u}\mathbf{f}}(\mathbf{z}, \mathbf{x}) &= \sum_{i=1}^m w_i \mathbb{I}\{\mathbf{z} \in \mathcal{Z}_i\} k_i(\mathbf{z}, \mathbf{x}), \\ k_{\mathbf{u}\mathbf{u}}(\mathbf{z}, \mathbf{z}') &= \sum_{i=1}^m \mathbb{I}\{\mathbf{z} \in \mathcal{Z}_i\} \mathbb{I}\{\mathbf{z}' \in \mathcal{Z}_i\} k_i(\mathbf{z}, \mathbf{z}'). \end{aligned} \quad (7)$$

From the posterior mean and covariance in Equation 1, we obtain the exact the formula in Equation 4.

One may concern about the approximate quality of MultiSVGP comparing to SVGP. Here, we argue that our approximate posterior can at least as good as SVGP. Let \hat{P} be the true posterior of GP, Q_{multi} be our variational approximation in Equation 4 and Q_{single} be the variational approximation of SVGP in Equation 3.

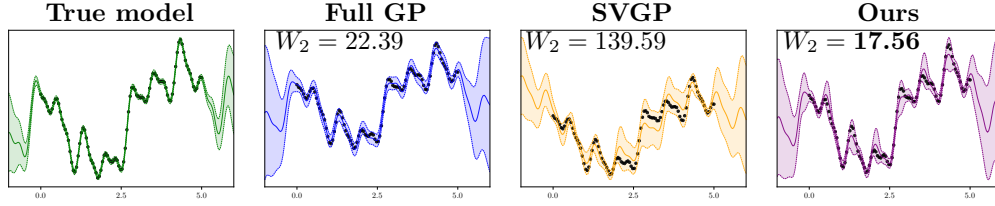


Figure 2: Comparison between the posterior of models. Here, W_2 measures the Wasserstein-2 distance between a model and the true model. The posterior obtained from our approach is close to the true model as well as the full GP model. SVGP model struggles to fit the data.

We only consider the case¹ $k(\mathbf{x}, \mathbf{x}') = k_1(\mathbf{x}, \mathbf{x}') + k_2(\mathbf{x}, \mathbf{x}')$. Let $\lambda_i, \lambda_i^{(1)}$, and $\lambda_i^{(2)}$ be the i -th operator eigenvalue w.r.t. k, k_1 , and k_2 . SVGP has M inducing points. MultiSVGP also has M inducing points in each inducing group. According to (Burt, Rasmussen, and Van Der Wilk 2019), the bound for the SVGP is

$$\text{KL}(Q_{\text{single}}||\hat{P}) \leq \frac{C_{\text{single}}}{2\sigma_n^2\delta} \left(1 + \frac{\|\mathbf{y}\|_2^2}{\sigma_n^2}\right),$$

with probability at least $1 - \delta$. Here, $C_{\text{single}} = \sum_{i=M+1}^{\infty} \lambda_i$, and σ_n^2 is a Gaussian noise variance. In the case of MultiSVGP, we bound the KL divergence by the following Proposition.

Proposition 1. *Given $k(\mathbf{x}, \mathbf{x}') = k_1(\mathbf{x}, \mathbf{x}') + k_2(\mathbf{x}, \mathbf{x}')$, with probability at least $1 - \delta$, we have*

$$\text{KL}(Q_{\text{multi}}||\hat{P}) \leq \frac{C_{\text{multi}}}{2\sigma_n^2\delta} \left(1 + \frac{\|\mathbf{y}\|_2^2}{\sigma_n^2}\right),$$

with $C_{\text{multi}} = N \sum_{i=1}^M (\lambda_i - \lambda_i^{(1)} - \lambda_i^{(2)}) + N \sum_{j=M+1}^{\infty} \lambda_j$.

Furthermore, it is true that $C_{\text{multi}} \leq C_{\text{single}}$, making the upper bound of $\text{KL}(Q_{\text{multi}}||\hat{P})$ is smaller or equal than the upper bound of $\text{KL}(Q_{\text{single}}||\hat{P})$.

Proof. We base on the result in (Burt, Rasmussen, and Van Der Wilk 2019) on the upper bound of the KL divergence and Equation 7 to derive $\text{KL}(Q_{\text{multi}}||\hat{P})$. We then use the property of the eigenvalues of sum matrices (Wielandt 1955; Tao 2010) where $\sum_{i=1}^M \lambda_i \leq \sum_{i=1}^M \lambda_i^{(1)} + \lambda_i^{(2)}$, to obtain the conclusion. The complete proof is placed in Appendix B. \square

This is considered a theoretical justification for the comparison in Figure 2.

5 Variational inference with shrinkage prior

With introducing the sparse vector \mathbf{w} in the kernel selection problem, we tackle the inference problem for this model with variational inference. That is, we consider the variational distribution which will be in the form of factorization between the approximate posterior distribution of GP latent variables and that of sparse vector \mathbf{w} . Let $q(\mathbf{w})$ be the variational distribution over \mathbf{w} .

¹without loss of generality, w_i is set to 1

The distribution $q(\mathbf{f}, \mathbf{U}, \mathbf{w}) = q(\mathbf{f}, \mathbf{U})q(\mathbf{w})$ approximates the true posterior. Following the approximate posterior construction of SVGP, $q(\mathbf{f}, \mathbf{U})$ is still in the similar form, $p(\mathbf{f}|\mathbf{U})q(\mathbf{U})$, with $q(\mathbf{U}) = \prod_{i=1}^m q(\mathbf{u}_i)$ and $q(\mathbf{u}_i)$ is parameterized by $\mathcal{N}(\mathbf{m}_i, \mathbf{S}_i)$. We maximize the evidence lower bound (ELBO)

$$\begin{aligned} \mathcal{L} &= \mathbb{E} \left[\log \frac{p(\mathbf{y}, \mathbf{f}, \mathbf{U}, \mathbf{w})}{q(\mathbf{f}, \mathbf{u}, \mathbf{w})} \right] \\ &= \mathbb{E}_{p(f(\cdot))} \left[\mathbb{E}_{q(\mathbf{w})} [\log p(\mathbf{y}|\mathbf{f}, \mathbf{w})] \right] \\ &\quad - \text{KL}(q(\mathbf{U})||p(\mathbf{U})) - \text{KL}(q(\mathbf{w})||p(\mathbf{w})). \end{aligned} \quad (8)$$

Note that $p(f(\cdot))$ is obtained after marginalizing all \mathbf{u}_i , can be done in the same manner with SVGP. The KL divergence w.r.t. \mathbf{U} is the sum of the individual KL divergence w.r.t. \mathbf{u}_i .

We describe the subroutine for variational inference w.r.t. \mathbf{w} , represented by Horseshoe variables, τ and λ . Due to the flat-tailed property of Half-Cauchy distribution, it is reparameterized by double inverse Gamma distributions (Wand et al. 2011; Ghosh, Yao, and Doshi-Velez 2019). That is, if $a \sim \mathcal{C}^+(b)$, this is equivalent to $a^2 \sim \text{IG}(1/2, \phi_a^{-1})$ and $\phi_a \sim \text{IG}(1/2, b^{-1})$. Now, we reintroduce the prior containing variables $\{\tau, \lambda, \phi_\tau, \phi_\lambda\}$ as

$$\begin{aligned} \tau^2 | \phi_\tau &\sim \text{IG}(1/2, \phi_\tau^{-1}), & \phi_\tau &\sim \text{IG}(1/2, A^{-1}), \\ \lambda_i^2 | \phi_{\lambda_i} &\sim \text{IG}(1/2, \phi_{\lambda_i}^{-1}), & \phi_{\lambda_i} &\sim \text{IG}(1/2, B^{-1}). \end{aligned}$$

We use the mean-field approach in which the variational distribution $q(\tau, \lambda, \phi_\tau, \phi_\lambda)$ is further factorized. Specifically, the variational distributions of τ and λ_i are chosen as log-normal distributions

$$\begin{aligned} q(\tau^2) &= \text{Lognormal}(\tau^2; \mu_\tau, \sigma_\tau^2), \\ q(\lambda_i^2) &= \text{Lognormal}(\lambda_i^2; \mu_{\lambda_i}, \sigma_{\lambda_i}^2). \end{aligned}$$

Whereas $q(\phi_\tau)$ and $q(\phi_{\lambda_i})$ remain inverse Gamma distribution.

As we exchange \mathbf{w} to $\{\tau, \lambda\}$, we have the expectation $\mathbb{E}_{q(\tau)q(\lambda)} [\log p(\mathbf{y}|\mathbf{f}, \tau, \lambda)]$ which is estimated by Monte Carlo integration. $q(\tau)$ and $q(\lambda)$ using reparameterization tricks for Log normal distributions (Kingma and Welling 2014). In particular, since the product $\tau^2 \lambda_i^2$ is also Log normal, it can be reparameterized by $\exp(\mu_\tau + \mu_{\lambda_i} + \varepsilon(\sigma_\tau + \sigma_{\lambda_i}))$ with $\varepsilon \sim \mathcal{N}(0, 1)$. We provide the detailed derivation of ELBO in Appendix C.

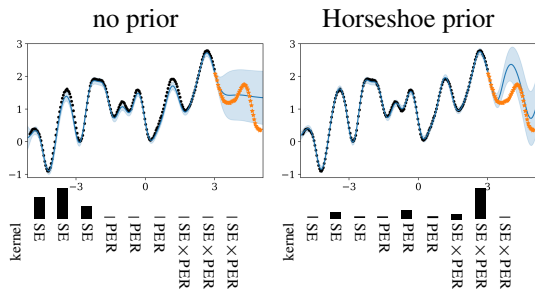


Figure 3: Behavior of Horseshoe prior in kernel selection. Both models predicts the test data (\star). The bar plots are the weights w_i corresponding to k_i .

Computational complexity Compared to SVGP using single inducing points, MultiSVGP takes $\mathcal{O}(m \max_i \{M_i^2\} b)$ at each optimization iteration with minibatch size b . Again, M_i is the number of inducing points \mathbf{u}_i .

6 Experimental Evaluations

In this section, we first set up choices for compositional kernels. We then justify how the Horseshoe assumption for kernel selection on synthetic data as well as time series data. Finally, we validate our model with regression and classification tasks. Our model is developed based on (Matthews et al. 2017).

6.1 Kernel function pool

Now, we present our approach in designing kernel structures for $\{k_i(\mathbf{x}, \mathbf{x}')\}$. A kernel function is constructed as a form of multiplicative kernel $\prod_{i=1}^{\alpha} \mathcal{B}_i$ where \mathcal{B}_i is a base kernel taking from SE, LIN, PER. In our experiment, the kernel pool is composed of all possible kernels having the multiplicative order up to 2. We allow duplication in kernels structure. The total number of kernels in the pool is 24. Each kernel in the pool remains interpretable and can be described by natural language explanations (Lloyd et al. 2014).

Hyperparameter initialization (Vanhatalo et al. 2013; Kim and Teh 2018) suggests two types of hyperparameter initialization: weak prior and strong prior. Unlike existing approaches requiring multiple restarts, we made sure our kernel pool covers both of them.

Behavior of Horseshoe prior To see how Horseshoe prior behaves for kernel selection problem, we created a synthetic data $(x_i, y_i)_{i=1}^{100}$ with $x_i \in [-5, 5]$ and y_i generated from kernel $\text{PER}_1 + \text{SE} \times \text{PER}_2$. We train our model and compare to the case where there is no prior on weights \mathbf{w} . Figure 3 shows our model spike at the relevant kernel structure ($\text{SE} \times \text{PER}$) while the model with no prior mistakenly assign weights for local variations (SE).

Small-sized 1D regression We verify our model on small-sized data sets: airline passenger volume, Mauna Loa CO₂ concentration. Figure 4 shows that our model can fit the data well. In the airline data set, the obtained kernel includes $\text{PER} \times \text{SE}$, LIN and SE while (Lloyd et al. 2014) reports $\text{LIN} + \text{PER} \times \text{SE} \times \text{LIN} + \text{SE}$ and a heteroscedastic

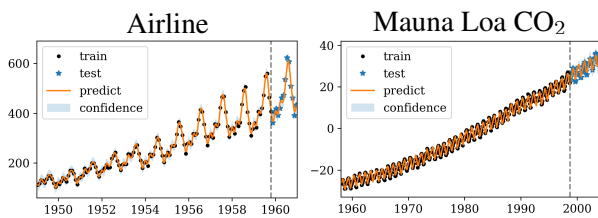


Figure 4: Extrapolation on time series data sets.

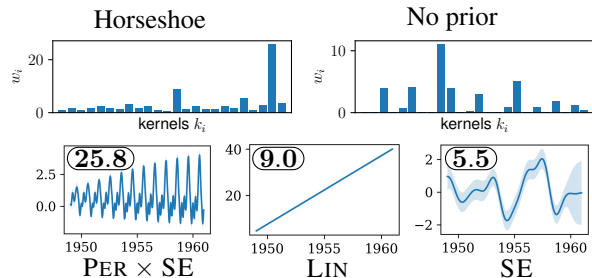


Figure 5: First row: the weights w_i in two cases. Second row: our kernel decomposition for airline data with three most significant components $\mathcal{GP}(\mu_i(\cdot), \Sigma_i(\cdot, \cdot))$. The weights w_i are showed at the upper-left corners.

noise. Also, in the mauna data set, the model can explain the trend and periodicity in data. Our model can reduce the running time to less than 0.5 hour comparing to 10 – 12 hours like (Duvenaud et al. 2013; Lloyd et al. 2014) or 2.5 – 4 hours like (Kim and Teh 2018). Figure 5 provides the visualization of weights \mathbf{w} found by our model comparing to the model without imposing any prior. This figure also gives the decomposition from $\mathcal{GP}(\sum w_i \mu_i(\cdot), w_i^2 \Sigma(\cdot, \cdot))$ corresponding with the most important components.

Medium-sized 1D regression We test our model on GEF-COM data set from the Global Energy Forecasting Competition (Tao Hong, Pierre Pinson, and Shu Fan 2014). The data set has $N = 38,070$ data points containing hourly records of energy load from January 2004 to June 2008. We randomly take 90% of the data set for training and held out 10% as test data. We compare our model with SVGP with no prior and SVGP with Softmax (Teng et al. 2020).

Figure 6 compares the predictive posteriors on the test set. It is clear that our model fits better, giving more accurate predictions as well as uncertainty estimation. The approach in (Teng et al. 2020) takes the second places. The inducing points are associated with complicated kernel function, not divided for each additive kernel. Therefore, the approximate capacity of this model is still more restricted than ours due to Proposition 1.

Our model found $\text{SE}_1 \times \text{PER}_1 + \text{SE}_2 + \text{SE}_3$ as the kernel structure for this data. This agrees with the kernel function in (Lloyd 2013) which is manually chosen. Also, our PER kernel has periodicity 1.001 days which also can describe the property that there are peaks in the morning and evening. This is aligned with the result reported in (Kim and Teh 2018).

Table 1: Extrapolation performance in UCI benchmarks. Results are aggregated from 10 independent runs.

	RMSE				Test log-likelihood			
	SVGP-SE	No prior	GP-NKN	Ours	SVGP-SE	No prior	GP-NKN	Ours
boston	7.30±0.21	7.24±0.27	5.53±0.49	5.41 ±0.10	-3.72±0.07	-3.72±0.10	-3.77±0.26	-3.24 ±0.11
concrete	9.64±0.14	8.70±1.05	6.44 ±0.19	7.39±0.42	-3.54±0.01	-3.45±0.08	-3.10 ±0.01	-3.33±0.06
energy	0.83±0.07	0.69±0.18	0.41±0.03	0.37 ±0.05	-1.11±0.03	-1.07±0.08	-0.54 ±0.04	-0.76±0.05
kin8nm	0.11±0.00	0.11±0.08	0.09 ±0.00	0.09 ±0.01	0.71±0.01	0.74±0.02	1.02 ±0.05	0.89±0.01
wine	0.62 ±0.00	0.62 ±0.01	0.67±0.01	0.63 ±0.01	-1.04±0.00	-1.04±0.00	-1.01 ±0.01	-1.04±0.01
yacht	1.45±0.10	1.22±0.44	0.46±0.05	0.36 ±0.05	-1.91±0.14	-1.67±0.46	-0.63 ±0.02	-0.83±0.12

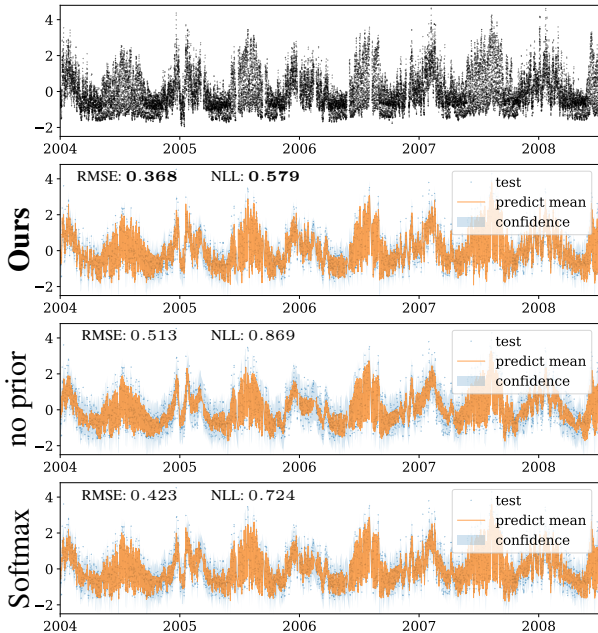


Figure 6: GEFCOM data set. First row is the plot of training data. The next rows are the predictive posterior at test points. Our model outperforms the alternatives in term of root mean square error (RMSE) and test negative log-likelihood (NLL).

Higher-dimension regression We conducted experiments on UCI data sets (Asuncion and Newman 2007) including boston, concrete, energy, kin8nm, wine and yacht (see Appendix D for detailed descriptions). We consider baseline models: GP-NKN (Sun et al. 2018), SVGP with no shrinkage prior over w (no prior), and SVGP with SE kernel (SVGP-SE). To justify the extrapolation performance, we projected data onto the principal component of data and sorted data according to the projection (Sun et al. 2018). From sorted indices, test data is taken from top $1/15$ and bottom $1/15$ of the data, the remaining is train data. We measure the root mean square error (RMSE) and test log-likelihood in each model.

Table 1 shows that our model has a competitive extrapolation capability comparing to GP-NKN. Roughly, our model has better performance in terms of RMSE for most of data sets, except concrete data set. In boston data set, our model performs well for the predictive log-likelihood. Still, GP-

NKN consistently outperforms others in this measurement. This is because this model is still considered as a full GP model retaining good uncertainty quantification while the remaining methods including ours are sparse GPs. However, GP-NKN takes significantly more time to train, e.g. in kin8nm. Although the model with no prior has a highly complex kernel, it fails to this extrapolation task. On the other hand, our model with shrinkage prior demonstrates the effect of regularization in kernel selection, resulting in better predictions.

Improving additive GPs (Duvenaud, Nickisch, and Rasmussen 2011) propose additive kernels for GPs to prevent the local property of kernel functions taking all input dimensions (Bengio, Delalleau, and Roux 2006). The additive kernel is the sum of lower-dimensional kernel functions which depend on a subset of input variables.

Suppose D is the dimension of data. Let $S_D = \{1, \dots, D\}$ be the index set of dimensions. We adopt this approach and consider the d -order additive kernel

$$k_d(\mathbf{x}, \mathbf{x}') = \sum_{\{i_1, \dots, i_d\} \subseteq S_D} w_{i_1 \dots i_d}^2 \prod_{i \in \{i_1, \dots, i_d\}} k(\mathbf{x}[i], \mathbf{x}'[i]).$$

Unlike (Duvenaud, Nickisch, and Rasmussen 2011) treating weights $w_{i_1 \dots i_d}$ equally in the same order d , we learn $\mathbf{w} = [w_{i_1 \dots i_d}]$ with our model. We conduct the experiment in three data sets: heart, liver, pima² (Duvenaud, Nickisch, and Rasmussen 2011) for classification task. The kernel type used here is SE kernel. The data sets are randomly split into training (90% of data) and test (10% of data) sets. We first run the model in (Duvenaud, Nickisch, and Rasmussen 2011) to obtain the most important order d . From d , we proceed learning $w_{i_1 \dots i_d}$. One limitation is that this setting is not scalable w.r.t D as the number of kernels, $\binom{D}{d}$, increases exponentially. Table 2 shows that our model can improve the accuracy of additive GPs by selecting appropriate kernels. On the other hand, the model without any prior even hurts the prediction. In the previous regression task, our model performs poorly in concrete data set since the 1-order additive kernels is the best fit for this data according to (Duvenaud, Nickisch, and Rasmussen 2011). We retrained the model and obtained an improved result with 6.90 ± 0.05 in RMSE, pushing the result closer to that of GP-NKN.

We provide an additional experiment comparing our model with the implicit kernel learning (Tompkins et al. 2019) in Appendix D.4.

²taken from <https://github.com/duvenaud/additive-gps>

Table 2: Classification error (in %) on three data sets.

	Additive GPs	No prior	Our model	#kernels $\binom{D}{d}$
Heart	18.15 \pm 4.56	16.00 \pm 1.41	14.00 \pm 2.11	$\binom{13}{1} = 13$
Liver	30.36 \pm 8.37	40.29 \pm 6.93	27.43 \pm 2.52	$\binom{6}{3} = 20$
Pima	23.99 \pm 3.46	29.87 \pm 3.72	20.52 \pm 1.65	$\binom{8}{6} = 28$

7 Conclusion

This paper presents a new approximate sparse GP targeting to improve compositional kernel learning for large-scaled data. Moreover, the paper presents a probabilistic kernel selection methods, showing satisfactory results in explaining time series as well as competitive extrapolation performance.

Acknowledgements

This work was supported by Institute of Information & communications Technology Planning & Evaluation (IITP) grant funded by the Korea government (MSIT) (No.2019-0-00075, Artificial Intelligence Graduate School Program (KAIST)). This work was partly done at VinAI Research. We would like to thank Juhyeong Kim for helping revising the paper, and anonymous reviewers for insightful feedback.

References

- Asuncion, A.; and Newman, D. 2007. UCI machine learning repository.
- Bauer, M.; van der Wilk, M.; and Rasmussen, C. E. 2016. Understanding Probabilistic Sparse Gaussian Process Approximations. In *NeurIPS*, 1533–1541.
- Bengio, Y.; Delalleau, O.; and Roux, N. L. 2006. The Curse of Highly Variable Functions for Local Kernel Machines. In *NeurIPS*, 107–114.
- Bhadra, A.; Datta, J.; Polson, N. G.; and Willard, B. 2019. Lasso Meets Horseshoe: A Survey. *Statist. Sci.* 34(3): 405–427.
- Bui, T. D.; Yan, J.; and Turner, R. E. 2017. A Unifying Framework for Gaussian Process Pseudo-Point Approximations using Power Expectation Propagation. *J. Mach. Learn. Res.* 18: 104:1–104:72.
- Burt, D.; Rasmussen, C. E.; and Van Der Wilk, M. 2019. Rates of Convergence for Sparse Variational Gaussian Process Regression. In *ICML*, 862–871.
- Carvalho, C. M.; Polson, N. G.; and Scott, J. G. 2009. Handling Sparsity via the Horseshoe. In *AISTATS*, 73–80.
- Duvenaud, D.; Lloyd, J. R.; Grosse, R.; Tenenbaum, J. B.; and Ghahramani, Z. 2013. Structure Discovery in Nonparametric Regression through Compositional Kernel Search. In *ICML*, 1166–1174.
- Duvenaud, D. K.; Nickisch, H.; and Rasmussen, C. E. 2011. Additive Gaussian Processes. In *NeurIPS*, 226–234.
- Ghahramani, Z. 2015. Probabilistic machine learning and artificial intelligence. *Nature* 521(7553): 452–459.
- Ghosh, S.; Yao, J.; and Doshi-Velez, F. 2019. Model Selection in Bayesian Neural Networks via Horseshoe Priors. *Journal of Machine Learning Research* 20(182): 1–46.
- Grosse, R.; Salakhutdinov, R.; Freeman, W. T.; and Tenenbaum, J. B. 2012. Exploiting compositionality to explore a large space of model structures. In *UAI*.
- Hensman, J.; Fusi, N.; and Lawrence, N. D. 2013. Gaussian Processes for Big Data. In *UAI*, 282–290. Arlington, Virginia, USA.
- Hwang, Y.; Tong, A.; and Choi, J. 2016. Automatic Construction of Nonparametric Relational Regression Models for Multiple Time Series. In *ICML*, 3030–3039.
- Kim, H.; and Teh, Y. W. 2018. Scaling up the Automatic Statistician: Scalable Structure Discovery using Gaussian Processes. In *AISTATS*, volume 84, 575–584.
- Kingma, D. P.; and Welling, M. 2014. Auto-Encoding Variational Bayes. In *ICLR*.
- Lázaro-Gredilla, M.; and Figueiras-Vidal, A. R. 2009. Inter-domain Gaussian Processes for Sparse Inference using Inducing Features. In *NeurIPS*, 1087–1095.
- Lloyd, J. R. 2013. GEFCom2012 Hierarchical Load Forecasting: Gradient boosting machines and Gaussian processes. *International Journal of Forecasting*.
- Lloyd, J. R.; Duvenaud, D.; Grosse, R.; Tenenbaum, J. B.; and Ghahramani, Z. 2014. Automatic Construction and Natural-Language Description of Nonparametric Regression Models. In *AAAI*, 1242–1250.
- Malkomes, G.; Schaff, C.; and Garnett, R. 2016. Bayesian optimization for automated model selection. In *NeurIPS*, 2900–2908.
- Matthews, A. G. d. G.; van der Wilk, M.; Nickson, T.; Fujii, K.; Boukouvalas, A.; León-Villagrà, P.; Ghahramani, Z.; and Hensman, J. 2017. GPflow: A Gaussian process library using TensorFlow. *J. Mach. Learn. Res.* 1–6.
- Mitchell, T. J.; and Beauchamp, J. J. 1988. Bayesian Variable Selection in Linear Regression. *Journal of the American Statistical Association* 83(404): 1023–1032.
- Neville, S. E.; Ormerod, J. T.; and Wand, M. P. 2014. Mean field variational Bayes for continuous sparse signal shrinkage: Pitfalls and remedies. *Electron. J. Statist.* 8(1): 1113–1151. doi:10.1214/14-EJS910.
- Paciorek, C. J.; and Schervish, M. J. 2003. Nonstationary Covariance Functions for Gaussian Process Regression. In *NeurIPS*, 273–280.
- Rasmussen, C. E.; and Ghahramani, Z. 2001. Occam’s Razor. In *NeurIPS*, 294–300.
- Rasmussen, C. E.; and Williams, C. K. I. 2005. *Gaussian Processes for Machine Learning*. The MIT Press.
- Saad, F. A.; Cusumano-Towner, M. F.; Schaechtle, U.; Rinaud, M. C.; and Mansinghka, V. K. 2019. Bayesian Synthesis of Probabilistic Programs for Automatic Data Modeling. *Proc. ACM Program. Lang.* 3(POPL): 37:1–37:32. doi:10.1145/3290350.

- Senanayake, R.; Tompkins, A.; and Ramos, F. 2018. Auto-morphing Kernels for Nonstationarity in Mapping Unstructured Environments. In *Proceedings of The 2nd Conference on Robot Learning*, 443–455.
- Snelson, E.; and Ghahramani, Z. 2006. Sparse Gaussian Processes using Pseudo-inputs. In *NeurIPS*, 1257–1264.
- Steinruecken, C.; Smith, E.; Janz, D.; Lloyd, J.; and Ghahramani, Z. 2019. The Automatic Statistician. In *Automated Machine Learning*.
- Sun, S.; Zhang, G.; Wang, C.; Zeng, W.; Li, J.; and Grosse, R. 2018. Differentiable Compositional Kernel Learning for Gaussian Processes. In *ICML*, 4828–4837.
- Tao, T. 2010. *254A, Notes 3a: Eigenvalues and sums of Hermitian matrices*. URL <https://terrytao.wordpress.com/2010/01/12/254a-notes-3a-eigenvalues-and-sums-of-hermitian-matrices/>.
- Tao Hong, Pierre Pinson, and Shu Fan. 2014. Global energy forecasting competition 2012.
- Teng, T.; Chen, J.; Zhang, Y.; and Low, B. K. H. 2020. Scalable Variational Bayesian Kernel Selection for Sparse Gaussian Process Regression. In *AAAI*, 5997–6004.
- Tibshirani, R. 1996. Regression Shrinkage and Selection via the Lasso. *Journal of the Royal Statistical Society (Series B)* 58: 267–288.
- Titsias, M. K.; and Lázaro-Gredilla, M. 2011. Spike and Slab Variational Inference for Multi-Task and Multiple Kernel Learning. In *NeurIPS*, 2339–2347.
- Tompkins, A.; Senanayake, R.; Morere, P.; and Ramos, F. 2019. Black Box Quantiles for Kernel Learning. In *Proceedings of Machine Learning Research*, 1427–1437.
- Tong, A.; and Choi, J. 2016. Automatic Generation of Probabilistic Programming from Time Series Data. *arXiv e-prints* arXiv:1607.00710.
- Tong, A.; and Choi, J. 2019. Discovering Latent Covariance Structures for Multiple Time Series. In *ICML*, 6285–6294.
- Vanhatalo, J.; Riihimäki, J.; Hartikainen, J.; Jylänki, P.; Tolvanen, V.; and Vehtari, A. 2013. GPstuff: Bayesian Modeling with Gaussian Processes. *J. Mach. Learn. Res.* 14(1).
- Wand, M. P.; Ormerod, J. T.; Padoan, S. A.; and Frühwirth, R. 2011. Mean Field Variational Bayes for Elaborate Distributions. *Bayesian Anal.* 6(4): 847–900. doi:10.1214/11-BA631. URL <https://doi.org/10.1214/11-BA631>.
- Wielandt, H. 1955. An extremum property of sums of eigenvalues .
- Williams, C. K. I.; and Seeger, M. 2001. Using the Nyström Method to Speed Up Kernel Machines. In *NeurIPS*, 682–688.
- Wilson, A. G.; and Adams, R. P. 2013. Gaussian Process Kernels for Pattern Discovery and Extrapolation. In *ICML*, 1067–1075.

A Kernels

Name	Kernel function $k(\mathbf{x}, \mathbf{x}')$
LIN	$(\mathbf{x} - \ell)^\top (\mathbf{x}' - \ell)$
SE	$\exp\left(-\frac{\ \mathbf{x} - \mathbf{x}'\ ^2}{2\ell^2}\right)$
PER	$\exp\left(-\frac{2 \sin^2(\pi \ \mathbf{x} - \mathbf{x}'\ /p)}{\ell^2}\right)$

Table 3: Base kernel functions. Note that we do not include the variance hyperparameters in these kernels since it is replaced by w_i .

B Proof for Proposition 1

Inter-domain covariance Before going to the proof of Proposition 1, we elucidate the derivation for Equation 7. We can represent our group inducing points as

$$u(\mathbf{z}) = \sum_{i=1}^m \int \mathbb{I}(\mathbf{z} \in \mathcal{Z}_i) \delta(\mathbf{z} - \mathbf{x}) g_i(\mathbf{x}) d\mathbf{x}.$$

Note that $\mathbb{I}(\mathbf{z} \in \mathcal{Z}_i)$ indicates the membership of inducing points in the group. If $\mathbf{z} \in \mathcal{Z}_i$, then $\mathbb{I}(\mathbf{z} \in \mathcal{Z}_i) = 1$. Otherwise, $\mathbb{I}(\mathbf{z} \in \mathcal{Z}_i) = 0$. The above representation means that only one member in the the group is selected with $\mathbb{I}(\mathbf{z} \in \mathcal{Z}_i) = 1$.

Following (Lázaro-Gredilla and Figueiras-Vidal 2009), we have the cross-covariance between $u(\mathbf{z})$ and $f(\mathbf{x})$ as

$$\begin{aligned}
 k_{\mathbf{uf}}(\mathbf{z}, \mathbf{x}) &= \mathbb{E}[u(\mathbf{z})f(\mathbf{x})] \\
 &= \mathbb{E} \left[\sum_{i=1}^m \int \mathbb{I}(\mathbf{z} \in \mathcal{Z}_i) \delta(\mathbf{z} - \mathbf{x}') g_i(\mathbf{x}') d\mathbf{x}' \sum_{j=1}^m w_j g_j(\mathbf{x}) \right] && \text{(by definition of } u \text{ and } f) \\
 &= \sum_{i=1}^m \sum_{j=1}^m w_j \int \mathbb{I}(\mathbf{z} \in \mathcal{Z}_i) \delta(\mathbf{z} - \mathbf{x}') \mathbb{E}[g_i(\mathbf{x}') g_j(\mathbf{x})] d\mathbf{x}' && \text{(swapping expectation inside sum and integration)} \\
 &= \sum_{i=1}^m w_i \mathbb{I}(\mathbf{z} \in \mathcal{Z}_i) \int \delta(\mathbf{z} - \mathbf{x}') \mathbb{E}[g_i(\mathbf{x}') g_i(\mathbf{x})] d\mathbf{x}' && \text{(since } \text{Cov}(g_i(\mathbf{x}'), g_j(\mathbf{x})) = 0 \text{ with } i \neq j) \\
 &= \sum_{i=1}^m w_i \mathbb{I}(\mathbf{z} \in \mathcal{Z}_i) \int \delta(\mathbf{z} - \mathbf{x}') k_i(\mathbf{x}', \mathbf{x}) d\mathbf{x}' && \text{(since } k_i(\mathbf{x}', \mathbf{x}) = \mathbb{E}[g_i(\mathbf{x}') g_i(\mathbf{x})]) \\
 &= \sum_{i=1}^m w_i \mathbb{I}(\mathbf{z} \in \mathcal{Z}_i) k_i(\mathbf{z}, \mathbf{x}). && \text{(integral with Diract delta function)}
 \end{aligned}$$

With similar steps, the covariance between u is obtained as

$$\begin{aligned}
k_{\mathbf{uu}}(\mathbf{z}, \mathbf{z}') &= \mathbb{E}[u(\mathbf{z})u(\mathbf{z}')] \\
&= \mathbb{E} \left[\sum_{i=1}^m \int \mathbb{I}(\mathbf{z} \in \mathcal{Z}_i) \delta(\mathbf{z} - \mathbf{x}) g_i(\mathbf{x}) d\mathbf{x} \sum_{j=1}^m \int \mathbb{I}(\mathbf{z}' \in \mathcal{Z}_j) \delta(\mathbf{z}' - \mathbf{x}') g_j(\mathbf{x}') d\mathbf{x}' \right] \\
&= \sum_{i=1}^m \sum_{j=1}^m \mathbb{I}(\mathbf{z} \in \mathcal{Z}_i) \mathbb{I}(\mathbf{z}' \in \mathcal{Z}_j) \int \int \delta(\mathbf{z} - \mathbf{x}) \delta(\mathbf{z}' - \mathbf{x}') \mathbb{E}[g_i(\mathbf{x}) g_j(\mathbf{x}')] d\mathbf{x} d\mathbf{x}' \\
&= \sum_{i=1}^m \mathbb{I}(\mathbf{z} \in \mathcal{Z}_i) \mathbb{I}(\mathbf{z}' \in \mathcal{Z}_i) \int \int \delta(\mathbf{z} - \mathbf{x}) \delta(\mathbf{z}' - \mathbf{x}') \mathbb{E}[g_i(\mathbf{x}) g_i(\mathbf{x}')] d\mathbf{x} d\mathbf{x}' \\
&= \sum_{i=1}^m \mathbb{I}(\mathbf{z} \in \mathcal{Z}_i) \mathbb{I}(\mathbf{z}' \in \mathcal{Z}_i) \int \int \delta(\mathbf{z} - \mathbf{x}) \delta(\mathbf{z}' - \mathbf{x}') k_i(\mathbf{x}, \mathbf{x}') d\mathbf{x} d\mathbf{x}' \\
&= \sum_{i=1}^m \mathbb{I}(\mathbf{z} \in \mathcal{Z}_i) \mathbb{I}(\mathbf{z}' \in \mathcal{Z}_i) k(\mathbf{z}, \mathbf{z}').
\end{aligned}$$

Proof of Proposition 1 In Lemma 1 of (Burt, Rasmussen, and Van Der Wilk 2019), we have

$$\text{KL}(Q_{\text{multi}} \|\hat{P}) \leq \frac{t}{2\sigma_n^2} \left(1 + \frac{\|\mathbf{y}\|_2^2}{\sigma_n^2 + t} \right).$$

Here $t = \text{trace}(\mathbf{K}_{\text{ff}} - \underbrace{\mathbf{K}_{\text{uf}}^\top \mathbf{K}_{\text{uu}}^{-1} \mathbf{K}_{\text{uf}}}_{\mathbf{Q}_{\text{ff}}})$. \mathbf{K}_{uf} is the cross-covariance matrix between inducing points \mathbf{U} and \mathbf{f} . \mathbf{K}_{uu} is the covariance matrix of \mathbf{U} .

Recall the covariance between \mathbf{U} and \mathbf{f} in Equation 7:

$$\begin{aligned}
k_{\text{uf}}(\mathbf{z}, \mathbf{x}) &= \mathbb{I}\{\mathbf{z} \in \mathcal{Z}_1\} k_1(\mathbf{z}, \mathbf{x}) + \mathbb{I}\{\mathbf{z} \in \mathcal{Z}_2\} k_2(\mathbf{z}, \mathbf{x}), \\
k_{\text{uu}}(\mathbf{z}, \mathbf{z}') &= \mathbb{I}\{\mathbf{z} \in \mathcal{Z}_1\} \mathbb{I}\{\mathbf{z}' \in \mathcal{Z}_1\} k_1(\mathbf{z}, \mathbf{z}') + \mathbb{I}\{\mathbf{z} \in \mathcal{Z}_2\} \mathbb{I}\{\mathbf{z}' \in \mathcal{Z}_2\} k_2(\mathbf{z}, \mathbf{z}'),
\end{aligned}$$

where $\mathbb{I}(\mathbf{z} \in \mathcal{Z}_i) = 1$ if $\mathbf{z} \in \mathcal{Z}_i$, otherwise it is equal 0. Then, we can compute

$$\begin{aligned}
\mathbf{K}_{\text{uf}} &= \text{Cov} \left(\begin{bmatrix} \mathbf{u}_1 \\ \mathbf{u}_2 \end{bmatrix}, \mathbf{f} \right) = \begin{bmatrix} \mathbf{K}_{\mathbf{u}_1 \mathbf{f}} \\ \mathbf{K}_{\mathbf{u}_2 \mathbf{f}} \end{bmatrix}, \\
\mathbf{K}_{\text{uu}} &= \text{Cov} \left(\begin{bmatrix} \mathbf{u}_1 \\ \mathbf{u}_2 \end{bmatrix}, \begin{bmatrix} \mathbf{u}_1 \\ \mathbf{u}_2 \end{bmatrix} \right) = \begin{bmatrix} \mathbf{K}_{\mathbf{u}_1 \mathbf{u}_1} & \mathbf{0} \\ \mathbf{0} & \mathbf{K}_{\mathbf{u}_2 \mathbf{u}_2} \end{bmatrix},
\end{aligned}$$

where $\mathbf{K}_{\mathbf{u}_1 \mathbf{f}}$ is the covariance matrix computed from $k_1(\cdot, \cdot)$ with $\mathbf{u}_1 \in \mathcal{Z}_1$, and $\mathbf{K}_{\mathbf{u}_1 \mathbf{u}_1}$ is the covariance matrix of $\mathbf{u}_1 \in \mathcal{Z}_1$ computed from $k_1(\cdot, \cdot)$. The same is applied to $\mathbf{K}_{\mathbf{u}_2 \mathbf{f}}$ and $\mathbf{K}_{\mathbf{u}_2 \mathbf{u}_2}$.

From $\mathbf{Q}_{\text{ff}} = \mathbf{K}_{\text{uf}}^\top \mathbf{K}_{\text{uu}}^{-1} \mathbf{K}_{\text{uf}}$ and block matrix multiplication, we can obtain

$$\mathbf{Q}_{\text{ff}} = \mathbf{K}_{\mathbf{u}_1 \mathbf{f}}^\top \mathbf{K}_{\mathbf{u}_1 \mathbf{u}_1}^{-1} \mathbf{K}_{\mathbf{u}_1 \mathbf{f}} + \mathbf{K}_{\mathbf{u}_2 \mathbf{f}}^\top \mathbf{K}_{\mathbf{u}_2 \mathbf{u}_2}^{-1} \mathbf{K}_{\mathbf{u}_2 \mathbf{f}}.$$

Let $\psi_i^{(1)}$ and $\psi_i^{(2)}$ be the eigenfunctions of the covariance operators w.r.t. $k_1(\mathbf{x}, \mathbf{x}')$ and $k_2(\mathbf{x}, \mathbf{x}')$. Following the similar steps in (Burt, Rasmussen, and Van Der Wilk 2019), we obtain the individual terms in \mathbf{Q}_{ff} under the eigenfeature representation as

$$\begin{aligned}
[\mathbf{K}_{\mathbf{f} \mathbf{u}_1} \mathbf{K}_{\mathbf{u}_1 \mathbf{u}_1}^{-1} \mathbf{K}_{\mathbf{u}_1 \mathbf{f}}]_{c,r} &= \sum_{i=1}^M \lambda_i^{(1)} \psi_i^{(1)}(\mathbf{x}_c) \psi_i^{(1)}(\mathbf{x}_r), \\
[\mathbf{K}_{\mathbf{f} \mathbf{u}_2} \mathbf{K}_{\mathbf{u}_2 \mathbf{u}_2}^{-1} \mathbf{K}_{\mathbf{u}_2 \mathbf{f}}]_{c,r} &= \sum_{i=1}^M \lambda_i^{(2)} \psi_i^{(2)}(\mathbf{x}_c) \psi_i^{(2)}(\mathbf{x}_r).
\end{aligned}$$

By Mercer's theorem and with $\psi_i(\mathbf{x})$ be the eigenfunctions of covariance operator for $k(\mathbf{x}, \mathbf{x}')$, we have

$$k(\mathbf{x}, \mathbf{x}') = \sum_{i=1}^{\infty} \lambda_i \psi_i(\mathbf{x}) \psi_i(\mathbf{x}').$$

Then the entry at (n, n) of $\mathbf{K}_{\text{ff}} - \mathbf{Q}_{\text{ff}}$ is

$$[\mathbf{K}_{\text{ff}} - \mathbf{Q}_{\text{ff}}]_{n,n} = \sum_{i=1}^{\infty} \lambda_i \psi_i^2(\mathbf{x}_n) - \left(\sum_{i=1}^M \lambda_i^{(1)} (\psi_i^{(1)})^2(\mathbf{x}_n) + \sum_{i=1}^M \lambda_i^{(2)} (\psi_i^{(2)})^2(\mathbf{x}_n) \right).$$

From this, we can have the expectation of t

$$\mathbb{E}_{\mathbf{x}}[t] = N \sum_{i=1}^{\infty} \lambda_i - N \left(\sum_{i=1}^M \lambda_i^{(1)} + \sum_{i=1}^M \lambda_i^{(2)} \right).$$

Here the eigenfunction terms are disappeared. Because $\mathbb{E}[\psi_i^2(\mathbf{x})] = \int \psi^2(\mathbf{x})p(\mathbf{x})d\mathbf{x} = 1$. Similarly, $\mathbb{E}[(\psi_i^{(1)})^2(\mathbf{x})] = \mathbb{E}[(\psi_i^{(2)})^2(\mathbf{x})] = 1$.

According to (Burt, Rasmussen, and Van Der Wilk 2019), we apply the Markov's inequality, we have, with probability at least $1 - \delta$,

$$\text{KL}(Q_{\text{multi}}||\hat{P}) \leq \frac{C_{\text{multi}}}{2\sigma^2\delta} \left(1 + \frac{\|\mathbf{y}\|_2^2}{\sigma_n^2} \right),$$

where $C_{\text{multi}} = N \sum_{i=1}^M (\lambda_i - \lambda_i^{(1)} - \lambda_i^{(2)}) + N \sum_{j=M+1}^{\infty} \lambda_j$.

Comparing to $C_{\text{single}} = \sum_{i=M+1}^{\infty} \lambda_i$, we use the result in (Wielandt 1955; Tao 2010) where

$$\sum_{i=1}^M \lambda_i \leq \sum_{i=1}^M \lambda_i^{(1)} + \sum_{i=1}^M \lambda_i^{(2)}.$$

We can conclude that the upper bound of $\text{KL}(Q_{\text{multi}}||\hat{P})$ is smaller than the upper bound of $\text{KL}(Q_{\text{single}}||\hat{P})$.

C Detail of variational inference

Prior Recall the prior over $\tau, \lambda, \phi_{\tau}, \phi_{\lambda}$ after reparameterization is

$$\begin{aligned} \tau^2 | \phi_{\tau} &\sim \text{IG}(\tau^2 | 1/2, \phi_{\tau}^{-1}), \\ \phi_{\tau} &\sim \text{IG}(\phi_{\tau} | 1/2, A^{-1}), \\ \lambda_i^2 | \phi_{\lambda_i} &\sim \text{IG}(\lambda_i^2 | 1/2, \phi_{\lambda_i}^{-1}), \quad i = 1 \dots m, \\ \phi_{\lambda_i} &\sim \text{IG}(\phi_{\lambda_i} | 1/2, B^{-1}). \end{aligned}$$

Variational distribution The variational distributions of τ and λ_i are in the form of log normal distribution.

$$\begin{aligned} q(\tau^2) &= \text{Lognormal}(\tau^2 | m_{\tau}, \sigma_{\tau}^2) \\ q(\lambda_i^2) &= \text{Lognormal}(\lambda_i^2 | m_{\lambda_i}, \sigma_{\lambda_i}^2), \quad i = 1 \dots m. \end{aligned}$$

On the other hand, the variational distributions of the auxiliary variables ϕ_{τ} and ϕ_{λ_i} remain as Inverse Gamma distributions

$$\begin{aligned} q(\phi_{\tau}) &= \text{IG}(\phi_{\tau} | s_{\tau}, r_{\tau}) \\ q(\phi_{\lambda_i}) &= \text{IG}(\phi_{\lambda_i} | s_{\lambda_i}, r_{\lambda_i}), \quad i = 1 \dots m. \end{aligned}$$

KL divergence As $\mathbf{w} = \{\tau, \lambda, \phi_{\tau}, \phi_{\lambda}\}$, the KL divergence $\text{KL}(q(\mathbf{w})||p(\mathbf{w}))$ becomes

$$\begin{aligned} &\text{KL} \left(q(\tau^2)q(\phi_{\tau}) \prod_i q(\lambda_i^2)q(\phi_{\lambda_i}) || p(\tau^2|\phi_{\tau})p(\phi_{\tau}) \prod_i p(\lambda_i|\phi_{\lambda_i})p(\lambda_i) \right) \\ &= H[q(\tau^2)] + H[q(\phi_{\tau})] + \sum_i H[q(\lambda_i^2)] + \sum_i H[q(\phi_{\lambda_i})] + \\ &\quad \mathbb{E}_{q(\tau^2)q(\phi_{\tau})}[\log p(\tau^2|\phi_{\tau})] + \mathbb{E}_{q(\phi_{\tau})}[\log p(\phi_{\tau})] + \sum_i \mathbb{E}_{q(\lambda_i)q(\phi_{\lambda_i})}[\log p(\lambda_i^2|\phi_{\lambda_i})] + \sum_i \mathbb{E}_{q(\phi_{\lambda_i})}[\log p(\phi_{\lambda_i})]. \end{aligned} \tag{9}$$

where $H[\cdot]$ denotes the entropy of a distribution.

Individual terms will be explained as following. The entropy terms will be computed as

$$\begin{aligned} H[q(\tau^2)] &= \mu_{\tau} + \frac{1}{2} \log(2\pi e \sigma_{\tau}^2), \\ H[q(\lambda_i^2)] &= \mu_{\lambda_i} + \frac{1}{2} \log(2\pi e \sigma_{\lambda_i}^2). \end{aligned}$$

The expectations of log prior can be derived as

$$\begin{aligned}
\mathbb{E}_{q(\tau^2)q(\phi_\tau)}[\log p(\tau^2|\phi_\tau)] &= \mathbb{E}_{q(\tau^2)q(\phi_\tau)}[\log \text{IG}(\tau^2|1/2, \phi_\tau^{-1})] \\
&= \mathbb{E}_{q(\tau^2)q(\phi_\tau)}\left[-\frac{1}{2}\log \phi_\tau - \log \Gamma(1/2) - \frac{3}{2}\log(\tau^2) - \frac{1}{\tau^2\phi_\tau}\right] \\
&= -\frac{1}{2}\mathbb{E}_{q(\phi_\tau)}[\log \phi_\tau] - \log \Gamma(1/2) - \frac{3}{2}\mathbb{E}_{q(\tau^2)}[\log(\tau^2)] - \mathbb{E}_{q(\tau^2)}[\tau^{-2}]\mathbb{E}_{\phi_\tau}[\phi_\tau^{-1}],
\end{aligned}$$

where the individual terms can be calculated as

$$\begin{aligned}
\mathbb{E}_{q(\phi_\tau)}[\log \phi_\tau] &= \log r_\tau - \psi(s_\tau), && \text{(Inverse Gamma distribution property)} \\
\mathbb{E}_{q(\phi_\tau)}[\phi_\tau^{-1}] &= \frac{s_\tau}{r_\tau}, && \text{(Inverse Gamma distribution property)} \\
\mathbb{E}_{q(\tau^2)}[\log(\tau^2)] &= \mu_\tau && \text{(compute from log normal distribution)} \\
\mathbb{E}_{q(\tau^2)}[\tau^{-2}] &= \exp(-\mu_\tau + \frac{1}{2}\sigma_\tau^2). && \text{(Log normal distribution property)}
\end{aligned}$$

Here, $\psi(\cdot)$ is the digamma function (is not the same with ψ in Section B). Similarly, we can obtain the expectation of log prior w.r.t to λ_i

$$\mathbb{E}_{q(\lambda_i^2)q(\phi_{\lambda_i})}[\log p(\lambda_i^2|\phi_{\lambda_i})] = -\frac{1}{2}\mathbb{E}_{q(\phi_{\lambda_i})}[\log \phi_{\lambda_i}] - \log \Gamma(1/2) - \frac{3}{2}\mathbb{E}_{q(\lambda_i^2)}[\log(\lambda_i^2)] - \mathbb{E}_{q(\lambda_i^2)}[\lambda_i^{-2}]\mathbb{E}_{\phi_{\lambda_i}}[\phi_{\lambda_i}^{-1}]$$

We intentionally do not write the explicit form of $H[q(\phi_\tau)]$, $H[q(\phi_{\lambda_i})]$, $\mathbb{E}_{q(\phi_\tau)}[\log p(\phi_\tau)]$ and $\mathbb{E}_{q(\phi_{\lambda_i})}[\log p(\phi_{\lambda_i})]$ because the variables ϕ_τ and ϕ_{λ_i} do not follow an optimization but are updated by the following.

Closed-form update for $q(\phi_\tau)$ and $q(\phi_{\lambda_i})$ Under the mean-field assumption on variational variables $\tau, \lambda, \phi_\tau, \phi_{\lambda_i}$, we can obtain the closed-form optimal solution w.r.t. the auxiliary variables $\phi_\tau, \phi_{\lambda_i}$ (Neville, Ormerod, and Wand 2014). That is, after each optimization step on other variables, we update $q(\phi_\tau)$ and $q(\phi_{\lambda_i})$ by

$$\begin{aligned}
q(\phi_\tau) &= \text{IG}(s_\tau = 1, r_\tau = \mathbb{E}[\tau^{-2}] + A^{-2}), \\
q(\phi_{\lambda_i}) &= \text{IG}(s_{\lambda_i} = 1, r_{\lambda_i} = \mathbb{E}[\lambda_i^{-2}] + B^{-2}).
\end{aligned} \tag{10}$$

Evidence lower bound Recap that the evidence lower bound is in the following form:

$$\mathcal{L} = \mathbb{E}_{p(f(\cdot))} [\mathbb{E}_{q(\tau, \lambda)}[\log p(\mathbf{y}|\mathbf{f}, \tau, \lambda)]] - \text{KL}(q(\mathbf{U})||p(\mathbf{U})) - \text{KL}(q(\tau, \lambda)||p(\tau, \lambda)). \tag{11}$$

Note that the expectation w.r.t τ, λ is estimated by Monte Carlo integration. During training, we draw one sample τ_S and λ_S by the reparameterization trick for the product $\tau_S \lambda_{i_S} = \exp(\mu_\tau + \mu_{\lambda_i} + \varepsilon(\sigma_\tau + \sigma_{\lambda_i}))$ where $\varepsilon \sim \mathcal{N}(0, 1)$.

The following algorithm describes our variational inference

Algorithm 1 Variational inference for MultiSVGP with Horseshoe prior

Require: Data $\mathcal{D} = \{\mathbf{X}, \mathbf{y}\}$, a set of kernel function $\{k_i(\mathbf{x}, \mathbf{x}')\}_{i=1}^m$

Initialize kernel hyperparameters, variational parameters $\{\mu_\tau, \sigma_\tau^2, \mu_{\lambda_i}, \sigma_{\lambda_i}^2, s_\tau, r_\tau, s_{\lambda_i}, r_{\lambda_i}\}$

for within a number of iterations **do**

 Sample a minibatch $(\mathbf{x}_b, \mathbf{y}_b)$

 Sample τ_S, λ_S with $\tau_S \lambda_{i_S} = \exp(\mu_\tau + \mu_{\lambda_i} + \varepsilon(\sigma_\tau + \sigma_{\lambda_i}))$ where $\varepsilon \sim \mathcal{N}(0, 1)$

 Compute $\mathbb{E}_{p(f(\mathbf{x}_b))} [\mathbb{E}_{q(\tau, \lambda)}[\log p(\mathbf{y}_b|\mathbf{f}, \tau, \lambda)]] \approx \mathbb{E}_{p(f(\mathbf{x}_b))} [\log p(\mathbf{y}_b|\mathbf{f}, \tau_S, \lambda_S)]$

 Compute $\text{KL}(q(\mathbf{U})||p(\mathbf{U}))$ as the sum of $\text{KL}(q(\mathbf{u}_i)||p(\mathbf{u}_i))$

 Compute $\text{KL}(q(\tau, \lambda)||p(\tau, \lambda))$ by Equation 9

 Compute ELBO \mathcal{L} based on Equation 11

 Perform an optimization step for ELBO \mathcal{L}

 Update $q(\phi_\tau)$ and $q(\phi_{\lambda_i})$ by Equation 10

end for

D Experiments

D.1 Description of regression data set

See Table 4

Table 4: Description of UCI data sets (Asuncion and Newman 2007)

Data set	# data N	Dimension D	Description
boston	506	13	Boston housing price
concrete	1030	8	Predict concrete compressive strength
energy	768	8	Predict energy efficiency for buildings
kin8nm	8192	8	Kinematics of an 8 link robot arm
wine	1599	22	Wine quality data set
yacht	308	7	Prediction of residuary resistance of sailing yachts

D.2 Description of classification task

See Table 5

Table 5: Description of heart, liver, pima data set

Data set	# data N	Dimension D	Description
heart	303	13	Predict the presence of heart disease
liver	345	6	Predict liver disorders
pima	768	8	Pima Indians Diabetes Database

D.3 Detailed experiment setup

Number of inducing points The number of inducing points used in airline and mauna loa data set is 100. In the remaining experiments, the number of inducing points is set to 200.

Horseshoe hyperparameter We choose the hyperparameters $A = 1, B = 1$ in Horseshoe prior. Varying these hyperparameters (between $[1, 3]$) does not affect much to the results since the maximum number of kernels is not big, about 24 (in regression task) to 28 (in classification task).

Experiment with softmax assumption (Teng et al. 2020) We follow the model candidates in (Teng et al. 2020) where there are 12 possible GP models to select. We cannot run the setting in which there are 144 candidate models described in (Teng et al. 2020). We implement this model by introducing a deterministic softmax weights to select model which is slightly different from (Teng et al. 2020). Yet, it still reflects the choice of model selection.

Table 6 includes the results considering the softmax assumption. Our model still outperforms the model with softmax in terms of both RMSE and test log-likelihood.

Table 6: Extrapolation performance in UCI benchmarks. Results are aggregated from 10 independent runs.

Data set	RMSE					Test log-likelihood				
	SVGP-SE	No prior	Softmax	GP-NKN	Ours	SVGP-SE	No prior	Softmax	GP-NKN	Ours
boston	7.30 \pm 0.21	7.24 \pm 0.27	6.90 \pm 0.34	5.53 \pm 0.49	5.41 \pm 0.10	-3.72 \pm 0.07	-3.72 \pm 0.10	-3.36 \pm 0.04	-3.77 \pm 0.26	-3.24 \pm 0.11
concrete	9.64 \pm 0.14	8.70 \pm 1.05	7.46 \pm 0.43	6.44 \pm 0.19	7.39 \pm 0.42	-3.54 \pm 0.01	-3.45 \pm 0.08	-3.39 \pm 0.04	-3.10 \pm 0.01	-3.33 \pm 0.06
energy	0.83 \pm 0.07	0.69 \pm 0.18	0.51 \pm 0.05	0.41 \pm 0.03	0.37 \pm 0.05	-1.11 \pm 0.03	-1.07 \pm 0.08	-0.89 \pm 0.04	-0.54 \pm 0.04	-0.76 \pm 0.05
kin8nm	0.11 \pm 0.00	0.11 \pm 0.08	0.13 \pm 0.00	0.09 \pm 0.00	0.09 \pm 0.01	0.71 \pm 0.01	0.74 \pm 0.02	0.60 \pm 0.01	1.02 \pm 0.05	0.89 \pm 0.01
wine	0.62 \pm 0.00	0.62 \pm 0.01	0.64 \pm 0.01	0.67 \pm 0.01	0.63 \pm 0.01	-1.04 \pm 0.00	-1.04 \pm 0.00	-1.07 \pm 0.01	-1.01 \pm 0.01	-1.04 \pm 0.01
yacht	1.45 \pm 0.10	1.22 \pm 0.44	1.81 \pm 0.53	0.46 \pm 0.05	0.36 \pm 0.05	-1.91 \pm 0.14	-1.67 \pm 0.46	-2.16 \pm 0.13	-0.63 \pm 0.02	-0.83 \pm 0.12

D.4 Compare to implicit kernel learning

This experiment compares between our proposed method and Black Box Quantiles (BBQ) (Tompkins et al. 2019). Data sets and experimental setups are followed (Tompkins et al. 2019). Table 7 shows the RMSE and MNLL of two methods. Our model shows competitive results in most data sets. On the other hand, in the intra-filling tasks on two data sets (pores and tread) BBQ is better. In these two data sets, our model comes with the second place according to the results of remaining alternative models in (Tompkins et al. 2019).

Table 7: Comparison to an implicit kernel learning approach

Data set	RMSE		MNL	
	BBQ	Ours	BBQ	Ours
CO2	0.068	0.157	-1.242	-1.649
Passenger	0.096	0.035	-0.610	-1.702
Concrete	0.124	0.086	-0.577	-0.742
Noise	0.138	0.074	-0.173	-0.922
Rubber	0.248	0.225	0.523	-0.010
Pores	0.256	0.574	0.335	0.898
Tread	0.114	0.123	-0.754	-0.644