
Revisiting the Sample Complexity of Sparse Spectrum Approximation of Gaussian Processes

Quang Minh Hoang

Department of Computer Science
Carnegie-Mellon University
Pittsburgh, PA 15213
qhoang@andrew.cmu.edu

Trong Nghia Hoang

MIT-IBM Watson AI Lab
IBM Research
Cambridge, MA 02142
nghiaht@ibm.com

Hai Pham

Language Technologies Institute
Carnegie-Mellon University
Pittsburgh, PA 15213
htpham@cs.cmu.edu

David P. Woodruff

Department of Computer Science
Carnegie-Mellon University
Pittsburgh, PA 15213
dwoodruf@cs.cmu.edu

Abstract

We introduce a new scalable approximation for Gaussian processes with provable guarantees which hold simultaneously over its entire parameter space. Our approximation is obtained from an improved sample complexity analysis for sparse spectrum Gaussian processes (SSGPs). In particular, our analysis shows that under a certain data disentangling condition, an SSGP’s prediction and model evidence (for training) can well-approximate those of a full GP with low sample complexity. We also develop a new auto-encoding algorithm that finds a latent space to disentangle latent input coordinates into well-separated clusters, which is amenable to our sample complexity analysis. We validate our proposed method on several benchmarks with promising results supporting our theoretical analysis.

1 Introduction

A Gaussian process (GP) [36] is a popular probabilistic kernel method for regression which has found applications across many scientific disciplines. Examples of such applications include meteorological forecasting, such as precipitation and sea-level pressure prediction [2]; sensing and monitoring of ocean and freshwater phenomena such as temperature and plankton bloom [7, 12]; traffic flow and mobility demand predictions over urban road networks [9, 10, 29]; flight delay predictions [15, 19, 20]; and persistent robotics tasks such as localization and filtering [43]. The broad applicability of GPs is in part due to their expressive Bayesian non-parametric nature which provides a closed-form prediction [36] in the form of a Gaussian distribution with formal measures of predictive uncertainty, such as entropy and mutual information criteria [27, 39, 44]. Such expressiveness makes GPs not only useful as predictive methods but also a go-to representation for active learning applications [24, 23, 27, 44] or Bayesian optimization [38, 45, 22, 16] that need to optimize for information gain while collecting training data.

Unfortunately, the expressive power of a GP comes at a cost of poor scalability (i.e., cubic time [36]) in the size of the training data (see Section 2.1 below), hence limiting its use to small datasets. This prevents GPs from being applied more broadly to modern settings with increasingly growing volumes of data [15, 19, 20]. To sidestep this limitation, a prevalent research trend is to impose sparse structural assumptions [33, 34] on the GP’s kernel matrix to reduce its multiplication and inversion cost, which comprises the main bulk of the training and inference complexity. This results in a broad family of

sparse Gaussian processes [17, 19, 28, 37, 40] that are not only computationally efficient but also amenable to various forms of parallelism [29] and distributed computation [13, 20, 21, 18], further increasing their efficiency.

Despite such advantages, the sparsification components at the core of these methods are heuristically designed and do not come with provable guarantees that explicitly characterize the interplay between approximation quality and computational complexity. This motivates us to develop a more robust, theoretically-grounded approximation scheme for GPs that is both provable and amenable to the many fast computation schemes mentioned above. More specifically, our contributions include:

1. An analysis of a new approximation scheme that generates a sparse spectrum approximation of a GP with provable bounds on its sample complexity, which practically becomes significantly small when the input data exhibits a certain clustering structure. Furthermore, the impact of the approximation on the resulting training and inference qualities is also formally analyzed (Section 3.1).
2. A data partitioning algorithm inspired from the above analysis, which learns a cluster embedding that reorients the input distribution while ensuring reconstructability of the original distribution (Section 3.3). We show that using sparse spectrum Gaussian processes (SSGP) in the embedded space requires fewer samples to achieve the same level of approximation quality. This also induces a linear feature map which enables efficient training and inference of GPs.
3. An empirical study on benchmarks that demonstrates the efficiency of the proposed method over existing works in terms of its approximation quality versus computational efficiency (Section 4).

2 Related Work

In this section we provide an overview of Gaussian processes (Section 2.1), followed by a succinct summary of their spectral representations (Section 2.2).

2.1 Gaussian Processes (GPs)

A Gaussian process [36] defines a probabilistic prior over a random function $f(x)$ defined by mean function $m(x) = 0^1$ and kernel function $k(x; x^0)$. These functions induce a marginal Gaussian prior over the evaluations $\mathbf{g} = [g(x_1) \dots g(x_n)]^\top$ on an arbitrary finite subset of inputs $\mathbf{x} = [x_1 \dots x_n]$. Let x be an unseen input whose corresponding output $g(x)$ we wish to predict. The Gaussian prior over $[g(x_1) \dots g(x_n) \ g(x)]^\top$ implies the following conditional distribution:

$$\mathbf{g} \mid \mathbf{g}(\mathbf{x}) \sim \mathcal{N}(\mathbf{g}; \mathbf{K}^{-1} \mathbf{g}(\mathbf{x}); \mathbf{k}(x; \mathbf{x}) \ \mathbf{K}^{-1} \mathbf{k}); \quad (1)$$

where $\mathbf{k} = [k(x; x_1) \dots k(x; x_n)]^\top$ and \mathbf{K} denotes the Gram matrix induced by $k(x; x^0)$ on $\mathbf{x} = [x_1 \dots x_n]$ for which $K_{ij} = k(x_i; x_j)$. For a noisy observation perturbed by Gaussian noise such that $\mathbf{y} \sim \mathcal{N}(\mathbf{g}(\mathbf{x}); \sigma^2)$, Eq. (1) above can be integrated with $\mathcal{N}(\mathbf{g}; \sigma^2)$ to yield:

$$\mathbf{g} \mid \mathbf{g}(\mathbf{x}) \sim \mathcal{N}(\mathbf{g}; (\mathbf{K} + \sigma^2 \mathbf{I})^{-1} \mathbf{y}; \mathbf{k}(x; \mathbf{x}) \ (\mathbf{K} + \sigma^2 \mathbf{I})^{-1} \mathbf{k}); \quad (2)$$

which explicitly forms the predictive distribution of a Gaussian process. The defining parameter of $k(x; x^0)$ (see Section 2.2) is crucial to the predictive performance and needs to be optimized via minimizing the negative log likelihood of:

$$\ell(\sigma) = \frac{1}{2} \log |\mathbf{K} + \sigma^2 \mathbf{I}| + \frac{1}{2} \mathbf{y}^\top (\mathbf{K} + \sigma^2 \mathbf{I})^{-1} \mathbf{y}; \quad (3)$$

where we now use the subscript to indicate that \mathbf{K} is a function of \mathbf{x} . In practice, both training and prediction incur $\mathcal{O}(n^3)$ processing cost, which prevents direct use of Gaussian processes on large datasets that might contain more than tens of thousands of training inputs.

2.2 Sparse Spectrum Gaussian Processes

Sparse spectrum Gaussian processes (SSGPs) [7, 28] exploit Theorem 1 below to re-express the Gaussian kernel $k(x; x^0) = \exp(-0.5(x - x^0)^\top \Sigma^{-1}(x - x^0))$ (where $\Sigma = \text{diag}[\frac{\sigma_1^2}{\lambda} \dots \frac{\sigma_d^2}{\lambda}]$) as an integration over a spectrum of cosine functions such that the integrating distribution (over the frequencies that parameterize these functions) is a multivariate Gaussian.

¹For simplicity, we assume a zero mean function since we can always re-center the training output around

Theorem 1 (Bochner Theorem) Let $k(x; x^0)$ denote a Gaussian kernel defined above and let $p(r) \sim \mathcal{N}(0; (4^{-2} \sigma^2)^{-1})$. It follows that:

$$k(x; x^0) = \mathbb{E}_r \left[\sum_{i=1}^m \cos(2r_i^> (x - x^0)_i) \right]; \quad (4)$$

where r is a d -dimensional random variable that parameterizes $\sum_{i=1}^m \cos(2r_i^> (x - x^0)_i)$. In practice, r is often referred to as the spectral frequency.

This allows us to approximate the original Gram matrix with a low-rank matrix K^0 constructed by a linear kernel $K^0(x; x^0) = \phi(x)^> \phi(x^0)$ with feature map $\phi(x) = [\phi_1(x) \dots \phi_{2m}(x)]^>$ comprising $2m$ basis trigonometric functions [7]. Each pair of odd- and even-index basis functions $\phi_{2i-1}(x) = \cos(2r_i^> x)$ and $\phi_{2i}(x) = \sin(2r_i^> x)$ is parameterized by the same sample of spectral parameter $r_i \sim p(r)$. For efficient computation m is often selected to be significantly smaller than d (i.e., the number of training examples). However, to guarantee that $K^0 \approx K$ with probability at least $1 - \delta$, m needs to be as large as $\Omega(n^2 = \frac{2}{\delta} \log(n/\delta))$ [31]², which makes the total prediction complexity much worse than that of a full GP.

Alternatively, one can use kernel sketching methods [32, 35] to generate feature maps that scale more favorably with the effective dimension of the kernel matrix, which empirically tends to be on the order of $\mathcal{O}(\log n)$. However, the pitfall of these methods is that without knowing the exact parameter configuration that underlies the data, they cannot sample from the true $p(r)$, which is necessary in their analyses. As such, existing random maps [35] that were generated based on this spectral construction often depend on a parameter initialization and their approximation quality is only guaranteed for that particular parameter setting instead of uniformly over the entire parameter space. This motivates us to revisit the sample complexity of SSGP from a setting which specifically searches for a reorientation of the input distribution such that the reoriented data exhibits a disentangled cluster structure. Such disentanglement provides a more sample-efficient bound as we show in our analysis in Section 3.1 below.

3 Provable Approximation of SSGPs with Improved Sample Complexity

We first show how a sparse spectrum Gaussian process (SSGP) can be approximated well with a provably low sample complexity. This is achieved by revisiting its sample complexity which, unlike prior work [3, 31, 35], explicitly characterizes and accounts for a certain set of data disentanglement conditions. Importantly, our new analysis (Section 3.1) yields practical bounds on both an SSGP's prediction and model evidence (Section 3.2) that hold with high probability uniformly over the entire parameter space. Furthermore, our analysis also inspires an encoding algorithm that finds a latent space to disentangle the encoded coordinates of data into well-separated clusters on which a sparse spectrum GP can approximate a GP provably well (Section 3.3). Our experiments show that such a latent space GP can be found for several real-world datasets (Section 4).

3.1 Practically Improved Sample Complexity for Sparse Spectrum Gaussian Processes

This section derives a new data-oriented feature map to approximate a Gaussian process parameterized with a Gaussian kernel. Unlike existing work which assumes knowledge of the true kernel parameters [3, 32, 35], our derivation remains oblivious to such parameters, and therefore holds universally over their entire candidate space. We assume that the GP prior of interest is of the form $\text{GP}(0; k(x; x^0))$ where $k(x; x^0)$ represents its Gaussian kernel in Section 2.2.

We give our analysis in three parts: (1) the spectral sampling scheme and a notion of approximation loss; (2) a set of practical data conditions which can be either observed from a raw data distribution or approximately imposed on the data via a certain embedding; (3) a theoretical analysis that delivers our key result that establishes an improved sample complexity when our data conditions are met.

3.1.1 Spectral Sampling Scheme and Spectral Loss

We show that $g(x)$ can be approximated by $g^0(x) = \sum_{i=1}^P g_i(x)$ with provable data-oriented guarantees where $g(x) \sim \text{GP}(0; (1 - \frac{1}{P})k_i(x; x^0))$. To achieve this, we first establish in Lemma 1

²See Theorem 6.28 in Chapter 6 of [31].

³In contrast, existing literature often generates bounds on either an SSGP's prediction or its model evidence (for training) for a single parameter configuration, which makes such an analysis only heuristic.

that the induced Gram matrix K of $k(x; x^0)$ on any dataset can be represented as an expectation over a space of induced Gram matrices $\{K_i\}_{i=1}^p$ produced by a corresponding space of random kernels $\{k_i(x; x^0)\}_{i=1}^p$.

Lemma 1. Let $k(x; x^0)$ and K denote a Gaussian kernel parameterized by (Section 2.2) and its induced Gram matrix on an arbitrary set of training inputs, respectively. There exists a space of random kernels $\{k_i(x; x^0)\}$ and a μ -independent distribution over K for which $K = E[K]$ where K denotes the induced Gram matrix on the same set of training inputs.

This follows directly from Theorem 1 above which states that $k(x; x^0) = E[\cos(2r^T(x - x^0))]$ where $r \sim N(0; (4^{-2} \sigma^2)^{-1})$. We can choose $\{k_i(x; x^0)\} = \{\cos(\langle r_i, x - x^0 \rangle)\}$ where $r_i \sim N(0; 1)$ which implies $k(x; x^0) = E[k_i(x; x^0)]$. Thus, $K = E[K]$ where the μ -independent parameter μ and K is the induced Gram matrix of. Leveraging the result of Lemma 1, a naïve analysis [1] using worst-case concentration bounds to derive a conservative estimate for a sufficient number of samples would require a prohibitively expensive sample complexity of $O(n^2 \log n)$.

Such analyses, however, often ignore the input distribution, which can be used to sample more selectively, thereby significantly reducing the sample complexity. This is demonstrated below in Theorem 2 which shows that when the input distribution exhibits a certain degree of compactness and separation (as defined in Conditions 1-3), we only require $O((\log^2 n - 2) \log \log(n - 2))$ sampled kernels $\{k_i\}_{i=1}^p$ indexed by $\{i\}_{i=1}^p$ to produce an average Gram matrix $K = \frac{1}{p} \sum_{i=1}^p K_i$ that is sufficiently close to K in spectral norm (see Definition 1) with probability at least $1 - \epsilon$.

Definition 1 (Spectral Closeness) Given $\epsilon > 0$, the symmetric matrices K and K^0 are ϵ -close if $\|K - K^0\|_2 \leq \epsilon$ where $\|K\|_2 = \lambda_{\max}(K - K^0)$ denotes the largest eigenvalue of $K - K^0$.

Thus, parameterizing the GP prior with K^0 instead of K allows us to derive an upper bound on the expected difference between their induced model evidence (for learning kernel parameters) and prediction losses (for testing) with respect to the same parameter setup (Theorem 3). Theorem 3 importantly exploits the fact that the bound in Theorem 2 holds universally over the entire space of parameters, which allows us to bound the prediction difference between the original and approximated GPs with respect to their own optimized parameters (that are not necessarily the same).

3.1.2 Practical Conditions on Data Distributions

We now outline key practical data conditions, which can be satisfied approximately via an encoding algorithm that transforms the input data into a latent space where such conditions are met. These conditions are necessary for deriving a practically improved sample complexity in Section 3.1.3.

Condition 1. For each parameter configuration $\theta = \text{diag}[\frac{2}{d_1}; \dots; \frac{2}{d_d}]$, there exists a mixture distribution $M(x; \mu = (\mu_1; \dots; \mu_b); \sigma = (\sigma_1; \dots; \sigma_b); c = (c_1; \dots; c_b))$ with at most $b = O(\log n)$ Gaussian components $N(x; c_i; \frac{2}{d_i} \sigma_i^{-1})$ over the data space with the mixing weights $\pi_i = 2^{-i}$ and variances $\sigma_i = O(\frac{1}{d_i})$ that generate the observed data in d -dimensional space.

Condition 2. The i^{th} Gaussian component as defined in Condition 1 above was used to generate the data points of the observed dataset. This can be substantiated easily with high probability given the above setup in Condition 1 that assigns selection probability $\pi_i = 2^{-i}$ to the i^{th} -component.

Condition 3. For each parameter configuration $\theta = \text{diag}[\frac{2}{d_1}; \dots; \frac{2}{d_d}]$, the mixture distribution of data in Condition 1 has sufficiently separated cluster centers. That is, for all i, j :

$$2^{i-2} (c_i - c_j)^2 > \frac{3}{2} \log \frac{2^a}{2^a - 1} \quad \text{where} \quad a = \frac{1}{\log 2} \log \frac{n^4}{n^4 - 4} \quad (5)$$

These conditions impose that the observed data can be separated into a number of clusters with exponentially growing sizes and concentration (see the small variances defined in Condition 1 and the imposed sizes of Condition 2). Intuitively, this means data points that belong to clusters with high concentration are responsible for kernel entries with high values whereas those in clusters with low concentration generate entries with low values. This is easy to see since high concentration reduces the distance between data points, thus increasing their kernel values and vice versa.

Furthermore, as imposed by Condition 2, clusters with high concentration also have denser population and induce kernel entries with high value. In addition, Condition 3 requires that clusters are well-separated, which implies that a large number of kernel entries are small and therefore can be

approximated cheaply. Together, these conditions form the foundations of our reduced complexity analysis for SSGP in Theorem 2. Interestingly, we show that such conditions also inspire the development of a probabilistic algorithm that finds an encoding of the input that (approximately) satisfies these conditions while preserving the statistical properties of the input (Section 3.3). This results in an improved sample complexity for SSGPs in practice (see Section 3.1.3).

3.1.3 Main Results

To understand the intuition why an improved sample complexity can be obtained, we note that when data is partitioned in clusters with different concentrations and sizes, the kernel entries are also partitioned into multiple value-bands with narrow width (i.e., low variance). Exploiting this, we can calibrate a significantly lower sample complexity for each band using concentration inequalities that improve with lower variance [11, 25].

Then, to combine these in-band sample complexities efficiently, we further exploit the data conditions in Section 3.1.2 to show that statistically, value bands with smaller width also tend to be populated more densely. This allows us to aggregate these in-band sample costs into an overall sample complexity with low cost. In practice, this also inspires an embedding algorithm (Section 3.3) that transforms the data in such a way that the distribution of their induced kernel entries will be denser in narrower bands, which is advantageous in our analysis.

Formally, let C be the set of all kernel entries indexed (u, v) in the Gram matrix K such that x_u and x_v belong to the same cluster and C^c be its complement. Also, let C be partitioned into b value-bands $C_i = \{(u, v) \mid f(u; v) \in [2^{-i}, 2^{-(i-1)})\}$ for $i \in [1 : b]$ and let $C_0 = \{(u, v) \mid f(u; v) \geq 2^{-b}\}$ be a band that is only populated by very large kernel entries. Theorem 2 below shows that we can construct a spectral approximation of K with arbitrarily high probability and low sample complexity.

Theorem 2. For any $\epsilon \in (0, 1]$ and $\delta \in (0, 1]$, if the training data has n data points and satisfies Conditions 1-3 above with respect to ϵ , then with probability at least $1 - \delta$, the approximation $K^0 = (1 - \epsilon) \sum_{i=1}^b K_i$ where $K_i \sim N(0, I)$ is ϵ -spectral close to K .

Proof Sketch. Our proof strategy is outlined below. The formal statements are spelled out in Appendix A.

First, with a proper choice of a clustering partition, the cross-cluster entries are guaranteed to be sufficiently small so as to be well-approximated by zero. We can then show with high probability that any kernel entry that corresponds to a pair of unique data points from the same cluster can be well-approximated with a sample complexity that scales favorably with the cluster's variance. In particular, we show that kernel values induced by data points generated by lower-variance clusters (see Condition 1) will have smaller approximation variances than those generated by data from higher-variance clusters and therefore require fewer samples to produce the same level of approximation.

Second, for certain configurations of mixture weights, Condition 2 asserts that the number of data points from each cluster is inversely proportional to the cluster variance, which implies that a small sample complexity is enough to approximate the majority of kernel entries. More specifically, Lemma 3 shows that when the input points are distributed into clusters with certain choices of variances σ_i^2 and at an inversely proportional rate ω_i^{-1} , then with high probability, over all clusters, the kernel entries (excluding those on the diagonal) associated with pairs in the cluster belong to their corresponding band.

Lemma 5 shows that for $\epsilon = O(\log^2 n / \log \log n)$, with probability $1 - \delta$, the total approximation error of all kernel entries in C will be at most ϵn , which implies with probability $1 - \delta$, the total approximation cost for items in C is at most ϵn . Next, Lemma 2 establishes that with the above data distribution, C accounts for ϵn entries while C^c accounts for $n^2 - \epsilon n$ entries, which needs to be approximated with error at most ϵn .

Finally, Lemma 4 shows that when the clusters are sufficiently well-separated (see Condition 3) any kernel value corresponding to an arbitrary data pair with points belonging to different clusters is guaranteed to be smaller than ϵn , which then guarantees a total error of at most ϵn when they are uniformly approximated with zero. Putting these together yields a total error with probability

⁴The intuition here is that kernel entries in narrower bands are cheaper (in term of sample cost) to approximate.

⁵The exact bounds defining the band can be found in Appendix A.

$\frac{1}{2}$, which implies K and K^0 are ϵ -spectrally close since $\|K - K^0\|_2 \leq \epsilon \|K\|_F$. Please see Appendix A for details.

3.2 Approximation Loss for Prediction and Model Evidence

In terms of prediction and model evidence approximation, our result holds simultaneously for all parameter configurations and is thus oblivious to the choice of parameters (see Theorem 3). While existing kernel sketch methods [32] generically achieve near-linear complexity for the approximate feature map, they often require knowledge of the parameters to construct the kernel approximations. In contrast, our result in Theorem 2 can be leveraged to bound the same prediction discrepancy when the original and approximated GPs use their own optimized parameter configurations, as shown in Theorem 4 below. To establish Theorem 4, however, we first establish an intermediate result that bounds the prediction and model evidence in the case when both the original and approximated GPs use the same parameter configurations.

Theorem 3. Let $\epsilon < 1$ be a user-specified confidence as defined previously in Theorem 2 and let \tilde{K} be an approximation to K for which $\|K - \tilde{K}\|_2 \leq \epsilon \|K\|_2$ with probability $1 - \epsilon$, uniformly over the entire parameter space. Then, with probability $1 - \epsilon$, the following hold:

$$E[g(x)] = (1 - \frac{\epsilon}{2}) E[g^0(x)] \quad \text{and} \quad V[g(x)] = (1 - \frac{\epsilon}{2}) V[g^0(x)] - \frac{\epsilon}{2} \quad (6)$$

where $\frac{\epsilon}{2}$ is the noise of the variance ($E[g^0(x)]$), and $g(x)$, $g^0(x)$ respectively denote the predictive distributions of the full GP and the approximated GP pertaining to an arbitrary test input

Proof. This follows directly from Lemma 7 and Lemma 8 in Appendix B. \square

Finally, Theorem 4 analyzes how close the approximated predictive mean is to the full GP predictive mean when both are evaluated at the optimizer of their respective training objective.

Theorem 4. Let $\epsilon < 1$ be a user-specified confidence as defined in Theorem 2. Let \tilde{K} denote an approximation to K for which $\|K - \tilde{K}\|_2 \leq \epsilon \|K\|_2$ with probability at least $1 - \epsilon$ uniformly over the entire parameter space. Let θ and θ^0 denote the optimal hyperparameters obtained by respectively minimizing the negative log likelihood of a full GP and the approximated GP. With probability $1 - \epsilon$, the following holds:

$$E[g^0(x; \theta^0)] = (1 - \epsilon) E[g(x; \theta)] + \epsilon C(\theta; \theta^0) \quad (7)$$

where $C(\theta; \theta^0)$ and $D(\theta; \theta^0)$ are constants with respect to θ ; θ^0 .

Proof. This follows immediately from Lemma 11 in Appendix B, which was built on the result of Theorem 3 above. This completes our loss analysis for SSGPs. \square

3.3 Optimizing Feature Map Complexity

We next present a practical probabilistic embedding algorithm that transforms the input data to meet the requirements of Condition 3. Our method is built on the rich literature of variational auto-encoders (VAEs) [6], which is a broad class of deep generative models that combine the rigor of Bayesian methods and rich parameterization of (deep) neural networks to discover (non-linear) low-dimensional embeddings of data while preserving their statistical properties. We first provide a short review on VAEs below, followed by an augmentation that aims to achieve the impositions in Conditions 1-3 above.

3.3.1 Variational Auto-Encoders (VAEs)

Let x be a random variable with density function $p(x)$. We want to learn a latent variable model $p(x; z) = p(z)p(z|x)$ that captures this generative process. The latent variable model comprises a fixed latent prior $p(z)$ and a parametric likelihood $p(x|z)$. To learn θ , we maximize the variational evidence lower-bound (ELBO) $L(\theta; \theta^0)$ of $\log p(x)$:

$$L(\theta; \theta^0) = E_z \int q(z|x) \log p(x|z) - \text{KL} [q(z|x) \parallel p(z)] \quad (8)$$

⁶[3, 32] achieves a complexity $\mathcal{O}(nm^2)$ where m scales with the effective dimension of the kernel matrix.

with respect to an arbitrary posterior surrogate $q(z|x)$ over the latent variable z . The ELBO is always a lower-bound $\log p(x)$ regardless of our choice of $q(z|x)$. This is due to the non-negativity of the KL divergence as seen in the first part of the above equation.

This can be viewed as a stochastic auto-encoder with $q(z|x)$ and $p(x|z)$ acting as the encoder and decoder, respectively. Here θ and ϕ characterize the neural network parameterization of these models. Their learning is enabled via a re-parameterization $q(z|x)$ that enables stochastic gradient ascent.

3.3.2 Re-conjugating Data via an Augmenting Variational Auto-Encoder

To augment the above VAE framework [26, 30] to account for the impositions in Condition 1 and 2, we ideally want to conjugate the parameterization of the above generative process to guarantee that the marginal posterior $p(z) = \int q(z|x)p(x)dx$ will manifest itself in the form of a mixture of Gaussians with the desired concentration and population densities as stated in Condition

However, it is often difficult to make such an imposition directly given that we typically have no prior knowledge of $p(x)$. We instead impose the desired structure on the latent $p(z)$ and then penalize the divergence between $q(z)$ and $p(z)$ while optimizing for the above ELBO in Eq(8). That is, we parameterize $p(z) = \sum_{i=1}^h N(z; c_i; \sigma_i^2) + \dots + \sum_{j=1}^b N(z; c_j; \sigma_j^2)$ where $\sigma_i^2 / \sigma_j^2 = 2^{i-2}$ (see Condition 2), which encodes the desired clustering structure. This is then reflected on the marginal posterior $q(z)$ via augmenting the above ELBO as,

$$L(x; \theta; \phi) = E_{z \sim q} \log p(x|z) - \text{KL}(q(z)||p(z)) - \text{KL}(q(z)||p(z)); \quad (9)$$

where the penalty term $\text{KL}(q(z)||p(z))$ serves as an incentive to encourage $q(z)$ to assume the same clustering structure $p(z)$. The parameter β can be manually set to adjust the strength of the incentive. To encourage separation among learned clusters (see Condition 3), we also add an extra penalty term to the above augmented ELBO,

$$L(x; \theta; \phi) = L(x; \theta; \phi) + \sum_{i \neq j} \text{KL}(N(z; c_i; \sigma_i^2) || N(z; c_j; \sigma_j^2)) \quad (10)$$

Once these clusters are learned, we can use the resulting encoding $q(z|x)$ to transform each training input x into its latent projection and subsequently train an SSGP on the latent space of (instead of training it on the original data space). Our previous analysis can then be applied to give the desired sample complexity. The empirical efficiency of the proposed method is demonstrated in Section 4 below. Note that the cost of training the embedding is linear in the number of data points and therefore does not noticeably affect our overall running time.

4 Experiments

Datasets. This section presents our empirical studies on two real datasets: (a) the ABALONE dataset [42] with 3000 data points which was used to train a model that predicts the age of abalone (number of rings on its shell) from physical measurements such as length, diameter, height, whole weight, shucked weight, viscera weight and shell weight; and (b) the GAS SENSOR dataset with 4 million data points [5, 6] which was used to train a model that predicts the CO concentration (ppm) from measurements of humidity, temperature, flow rate, heater voltage and the resistant measures of gas sensors.

In both settings, we compare our revised SSGP method with the traditional SSGP on both datasets to demonstrate its sample efficiency. In particular, our SSGP method is applied on the embedded space of data which was generated and conjugated using the auto-encoding method in Section 3.3.2 to approximately meet the aforementioned Condition 3.

The detailed parameterization of our entire algorithm is provided in Appendix C. The traditional SSGP method on the other hand was applied directly to the data space. The prediction root-mean-square-error (RMSE) achieved by each method is reported at different sample complexities in Figure 1 below. All reported performances were averaged over independent runs on a computing server with a Tesla K40 GPU with 12GB RAM.

Results and Discussions. It can be observed from the results that at all levels of sample complexity, the revised SSGP achieves substantially better performance than its vanilla SSGP counterpart. This

⁷Our experimental code is released at https://github.com/hqminh/gp_sketch_nips.

(a) (b) (c)

Figure 1: Graphs of performance comparisons between our revised SSGP and the traditional SSGP on the ABALONE dataset [42] at varying sample complexities (see Theorem 2, 32 and 64).

is expected since our revised SSGP is guaranteed to require many fewer samples than the vanilla SSGP when the data is recon gured to exhibit a certain clustering structure (see Condition 1 and Theorem 2). As such, when both are set to operate at the same level of sample complexity, one would expect the revised SSGP to achieve better performance since SSGP generally performs better when its sample complexity is set closer to the required threshold. On the larger GAS SENSOR dataset (which contains approximately 4M data points), we also observe the same phenomenon from the performance comparison graph as shown in Figure 2a below: A vanilla SSGP needs to increase its number of samples to marginally improve its predictive performance while our revisited SSGP is able to outperform the former with the least number of samples (6).

(a) (b) (c)

Figure 2: Graphs of (a) performance comparison between our revisited SSGP's (with sample complexity $p = 16$) and the vanilla SSGP's (with sample complexity 16; 32; 64) on the GAS SENSOR dataset [4]; and visualizations of (b) original and (c) recon gured data distributions of GAS SENSOR data on \mathbb{R}^2 -dimensional latent space generated by our auto-encoding algorithm in Section 3.3.

Furthermore, a closer look into the data distribution (visualized on a 2D space in Fig. 2b) and the data recon gured data distribution (visualized on a 2D space in Fig. 2c) also corroborates our hypothesis earlier that a well-separated data partition with high in-cluster concentration (in the form of a mixture of clusters – see Condition 1) can be found (by our embedding algorithm in Section 3.3) to recon gure our data distribution to (approximately) meet the necessary technical conditions that enable our sample-complexity enhancement analysis (see Section 3.1). Due to limited space, interested readers are referred to Appendix D for more detailed empirical studies and demonstrations.

5 Conclusion

We present a new method and analysis for approximating Gaussian processes. We obtain provable guarantees for both training and inference, which are the first to hold simultaneously over the entire space of kernel parameters. Our results complement existing work in kernel approximation that often assumes knowledge of its defining parameters. Our results also reveal important (practical) insights that allow us to develop an algorithmic handle on the tradeoff between approximation quality and sample complexity, which is achieved via finding an embedding that disentangles the latent coordinates of data. Our empirical results show for many datasets, such a disentangled embedding space can be found, which leads to a significantly reduced sample complexity of SSGP.

6 Statement of Broader Impact

Our work focuses on approximating Gaussian processes using a mixture of practical methods and theoretical analysis to reconstruct data in ways that reduce their approximation complexity. As such, it could have significant broader impact by allowing users to more accurately solve practical problems such as the ones discussed in our introduction, while still providing concrete theoretical guarantees. While applications of our work to real data could result in ethical considerations, this is an indirect (and unpredictable) side-effect of our work. Our experimental work uses publicly available datasets to evaluate the performance of our algorithms; no ethical considerations are raised.

7 Acknowledgement

T. N. Hoang is supported by the MIT-IBM Watson AI Lab, IBM Research. Q. M. Hoang is supported by the Gordon and Betty Moore Foundation's Data-Driven Discovery Initiative through Grant GBMF4554, by the US National Science Foundation (DBI-1937540), by the US National Institutes of Health (R01GM122935), and by the generosity of Eric and Wendy Schmidt by recommendation of the Schmidt Futures program. D. Woodruff is supported by National Institute of Health grant 5R01HG 10798-2, Office of Naval Research grant N00014-18-1-2562, and a Simons Investigator Award.

References

- [1] Rakshit Allamraju and Girish Chowdhary. Communication efficient decentralized gaussian process fusion for multi-uas path planning. *American Control Conference*, pages 4442–4447, 05 2017.
- [2] T. J. Ansell et al. Daily mean sea level pressure reconstructions for the European-North Atlantic region for the period 1850-2003. *J. Climate*, 19(12):2717–2742, 2006.
- [3] H. Avron, M. Kapralov, C. Musco, C. Musco, A. Velingker, and A. Zandieh. Random fourier features for kernel ridge regression: Approximation bounds and statistical guarantees. In *ICML*, pages 253–262, 2017.
- [4] Javier Burgues. Gas Sensor Array Temperature Modulation Dataset, <https://archive.ics.uci.edu/ml/machine-learning-databases/00487/>.
- [5] Javier Burgués, Juan Manuel Jiménez-Soto, and Santiago Marco. Estimation of the limit of detection in semiconductor gas sensors through linearized calibration methods. *Analitica chimica acta*, 1013:13–25, 2018.
- [6] Javier Burgués and Santiago Marco. Multivariate estimation of the limit of detection by orthogonal partial least squares in temperature-modulated mox sensors. *Analitica chimica acta*, 1019:49–64, 2018.
- [7] N. Cao, K. H. Low, and J. M. Dolan. Multi-robot informative path planning for active sensing of environmental phenomena: A tale of two algorithms. *Proc. AAMAS*, pages 7–14, 2013.
- [8] J. Chen, N. Cao, K. H. Low, R. Ouyang, C. K.-Y. Tan, and P. Jaillet. Parallel Gaussian process regression with low-rank covariance matrix approximations. *Proc. UAI*, pages 152–161, 2013.
- [9] J. Chen, K. H. Low, C. K.-Y. Tan, A. Oran, P. Jaillet, J. M. Dolan, and G. S. Sukhatme. Decentralized data fusion and active sensing with mobile sensors for modeling and predicting spatiotemporal traffic phenomena. *Proc. UAI*, pages 163–173, 2012.
- [10] Jie Chen, Kian Hsiang Low, and Colin Tan. Gaussian process-based decentralized data fusion and active sensing for mobility-on-demand systems. *Robotics: Science and Systems*, 2013.
- [11] H. Chernoff. A measure of asymptotic efficiency for tests of hypothesis based on the sum of observations. *Annals of Mathematical Statistics*, 23:493–509, 1952.

- [12] J. M. Dolan, G. Podnar, S. Stancliff, K. H. Low, A. Elfes, J. Higinbotham, J. C. Hosler, T. A. Moisan, and J. Moisan. Cooperative aquatic sensing using the telesupervised adaptive ocean sensor net. In *Proc. SPIE Conference on Remote Sensing of the Ocean, Sea Ice, and Large Water Regions*, volume 7473, 2009.
- [13] Y. Gal, M. van der Wilk, and C. Rasmussen. Distributed variational inference in sparse Gaussian process regression and latent variable models. In *Proc. NIPS*, 2014.
- [14] Yarin Gal and Richard Turner. Improving the gaussian process sparse spectrum approximation by representing uncertainty in frequency inputs. 2015.
- [15] J. Hensman, N. Fusi, and N. D. Lawrence. Gaussian processes for big data. *ICML*, pages 282–290, 2013.
- [16] M. Hoang and C. Kingsford. Optimizing dynamic structures with bayesian generative search. In *International Conference on Machine Learning*, 2020.
- [17] Q. M. Hoang, T. N. Hoang, and K. H. Low. A generalized stochastic variational Bayesian hyperparameter learning framework for sparse spectrum Gaussian process regression. In *AAAI*, pages 2007–2014, 2017.
- [18] Q. M. Hoang, T. N. Hoang, K. H. Low, and C. Kingsford. Collective model fusion for multiple black-box experts. In *Proc. ICML*, 2019.
- [19] T. N. Hoang, Q. M. Hoang, and K. H. Low. A unifying framework of anytime sparse Gaussian process regression models with stochastic variational inference for big data. In *Proc. ICML*, pages 569–578, 2015.
- [20] T. N. Hoang, Q. M. Hoang, and K. H. Low. A distributed variational inference framework for unifying parallel sparse Gaussian process regression models. In *Proc. ICML*, pages 382–391, 2016.
- [21] T. N. Hoang, Q. M. Hoang, K. H. Low, and J. P. How. Collective online learning of Gaussian processes in massive multi-agent systems. In *Proc. AAAI*, 2019.
- [22] T. N. Hoang, Q. M. Hoang, O. Ruofei, and K. H. Low. Decentralized high-dimensional bayesian optimization with factor graphs. In *Proc. AAAI*, 2018.
- [23] T. N. Hoang, K. H. Low, P. Jaillet, and M. Kankanhalli. Nonmyopic Bayes-optimal active learning of Gaussian processes. In *Proc. ICML*, pages 739–747, 2014.
- [24] T. N. Hoang, K. H. Low, P. Jaillet, and M. S. Kankanhalli. Active learning is planning: Nonmyopic -Bayes-optimal active learning of Gaussian processes. In *ECML-PKDD Nectar Track*, pages 494–498, 2014.
- [25] W. Hoeffding. Probability inequalities for the sum of bounded random variables. *Journal of the American Statistical Association*, 58:13–30, 1963.
- [26] D. Kingma and M. Welling. Auto-Encoding Variational Bayes. *Proc. ICLR*, 2013.
- [27] A. Krause and C. Guestrin. Nonmyopic active learning of Gaussian processes: An exploration-exploitation approach. In *Proc. ICML*, pages 449–456, 2007.
- [28] M. Lázaro-Gredilla, J. Quiñero-Candela, C. E. Rasmussen, and A. R. Figueiras-Vidal. Sparse spectrum Gaussian process regression. *Journal of Machine Learning Research*, pages 1865–1881, 2010.
- [29] K. H. Low, J. Yu, J. Chen, and P. Jaillet. Parallel Gaussian process regression for big data: Low-rank representation meets Markov approximation. In *Proc. AAAI*, pages 2821–2827, 2015.
- [30] Emile Mathieu, Tom Rainforth, Siddharth Narayanaswamy, and Yee Whye Teh. Disentangling disentanglement in variational autoencoders. *ICML*, 2019.
- [31] Mehryar Mohri, Afshin Rostamizadeh, and Ameet Talwalkar. *Foundations of Machine Learning*. MIT press, 2018.

- [32] C. Musco and C. Musco. Recursive sampling for the nystrom method. *Proc. NIPS*, 2016.
- [33] J. Quiñero-Candela and C. E. Rasmussen. A unifying view of sparse approximate Gaussian process regression. *Journal of Machine Learning Research*, 6:1939–1959, 2005.
- [34] J. Quiñero-Candela, C. E. Rasmussen, and C. K. I. Williams. Approximation methods for gaussian process regression. *Large-Scale Kernel Machines*, pages 203–223, 2007.
- [35] A. Rahimi and B. Recht. Random features for large-scale kernel machines. *Proc. NIPS*, 2007.
- [36] C. E. Rasmussen and C. K. I. Williams. *Gaussian Processes for Machine Learning*. MIT Press, 2006.
- [37] M. Seeger, C. K. I. Williams, and N. D. Lawrence. Fast forward selection to speed up sparse Gaussian process regression. *Proc. AISTATS*, 2003.
- [38] J. Snoek, L. Hugo, and R. P. Adams. Practical Bayesian optimization of machine learning algorithms. In *Proc. NIPS*, pages 2960–2968, 2012.
- [39] N. Srinivas, A. Krause, S. Kakade, and M. Seeger. Gaussian process optimization in the bandit setting: No regret and experimental design. *Proc. ICML*, pages 1015–1022, 2010.
- [40] M. K. Titsias. Variational learning of inducing variables in sparse Gaussian processes. In *AISTATS*, 2009.
- [41] L. J. P. van der Maaten and G. E. Hinton. Visualizing high-dimensional data using t-SNE. *Journal of Machine Learning Research*, 9:2579–2605, 2008.
- [42] Sam Waugh. Abalone Dataset, howpublished <https://archive.ics.uci.edu/ml/machine-learning-databases/abalone/>.
- [43] N. Xu, K. H. Low, J. Chen, K. K. Lim, and E. B. Özgül. GP-Localize: Persistent mobile robot localization using online sparse Gaussian process observation models. In *Proc. AAAI*, pages 2585–2592, 2014.
- [44] Y. Zhang, T. N. Hoang, K. H. Low, and M. Kankanhalli. Near-optimal active learning of multi-output Gaussian processes. *Proc. AAAI*, pages 2351–2357, 2016.
- [45] Y. Zhang, T. N. Hoang, K. H. Low, and M. Kankanhalli. Information-based multi-fidelity bayesian optimization. In *NIPS Workshop Bayesian Optimization*, 2017.

A Intermediate Results for Theorem 2

Let $(x_u; x_v)$, $jK(x_u; x_v) = K^Q(x_u; x_v)$ where $K^Q(x_u; x_v) = (1-p) \prod_{i=1}^p K_i(x_u; x_v)$, and where $i \sim N(0,1)$ as defined in Lemma 1 above. We will first measure the approximation loss across different value-bands of $K(x_u; x_v)$, thereby deriving tight sample bounds for each band. Combining these with the union bound allows us to establish a much cheaper overall sample complexity as compared to the naïve $O(n^2 \log n)$ bound.

Lemma 2. Suppose the data distribution follows Conditions 1 and 2 above. Let (x_u) denote the cluster index of each data point x_u . Let $C_i = \{u; v \mid c(x_u) = c(x_v) = i\}$ and $C^0 = \{u; v \mid c(x_u) \neq c(x_v)\}$ denote the sets of in-cluster and out-cluster kernel entries, respectively, where $|C_i| = \frac{n^2}{4}$ and $|C^0| = \frac{3n^2}{4}$.

Proof. By Condition 2, since data points are scattered across clusters and each cluster has 2^{i-2} points, it follows that:

$$\begin{aligned} n &= \sum_{i=1}^b 2^{i-2} = \frac{2^{b+1} - 2}{2} \\ |C_i| &= \sum_{i=1}^b 2^{i-2} = 2^{b+1} - 1 = n \frac{2^{b+1} - 2}{2^{b+1} - 1} \approx \frac{n}{2} \\ |C^0| &= n^2 - |C_i| \approx \frac{3n^2}{4} \end{aligned} \quad (11)$$

This also implies that $b = O(\log n)$, which is consistent with Condition 1 above. \square

Lemma 3. Let $C_i = \{u; v \mid c(x_u) = c(x_v) = i\}$ for $i \in [1 : : b]$. Then with probability at least $1 - O(\exp(-\log n \frac{1}{d}))$, the following holds for all $(u; v) \in C_i$ for which $u \neq v$:

$$1 - \frac{1}{2^{a+i-1}} \leq K(x_u; x_v) < 1 + \frac{1}{2^{a+i-1}} \quad \text{where } a = \frac{1}{\log 2} \log \frac{n^4}{n^4 - 4} \quad (12)$$

Proof. If x_u and x_v are both generated from component i of the data distribution as defined in Condition 1, it follows that $(x_u - x_v) \sim N(c_i; \frac{2}{i}I)$. Therefore, by standard chi-squared tail bounds, with probability at least $1 - 2e^{-t}$, we have:

$$\frac{(x_u - x_v)^2}{2} = \frac{2}{i} \chi^2_d \leq \frac{2}{i} (d + \sqrt{2dt}) \quad (13)$$

where d is the data dimension. Using this, we can then figure out a setting for t such that $K(x_u; x_v)$ follows the above condition in Eq. 12. In particular, set

$$L(i) = \log \frac{2^{a+i-1}}{2^{a+i-1} - 1} \quad \text{and} \quad U(i) = \log \frac{2^{a+i-1} + 1}{2^{a+i-1} - 1} \quad (14)$$

We can then choose:

$$\frac{2}{i} = \frac{1}{4d} (U(i) + L(i)) \quad \text{and} \quad t = \frac{p}{d} \frac{U(i) - L(i)}{U(i) + L(i)} \leq \frac{p}{d} \quad (15)$$

so that by plugging these choices in Eq. 13 above, we have with probability at least $1 - 2e^{-t}$:

$$\begin{aligned} \frac{(x_u - x_v)^2}{2} &\leq \frac{1}{2} \log \frac{2^{a+i-1}}{2^{a+i-1} - 1} + \frac{1}{2} \log \frac{2^{a+i-1} + 1}{2^{a+i-1} - 1} \\ &\leq \frac{1}{2^{a+i-1}} \leq \frac{1}{2^{a+i-1}} \quad \# \\ |K(x_u; x_v) - 1| &\leq \frac{1}{2^{a+i-1}} \leq \frac{1}{2^{a+i-1}} \quad (16) \end{aligned}$$

Now, note that for any i for which $\frac{2}{i} \leq 4e^{-\frac{p}{d}}$ or $\frac{2}{i} \leq 4e^{-\frac{p}{d}} - 8i$, we have $\frac{2}{i} \leq 2e^{-t}$ since $t = \frac{p}{d}$. This also means $O(\exp(-\log n \frac{1}{d}))$ since $b = O(\log n)$.

That is, Eq.(16) and hence, Eq.(12), hold with probability at least $1 - 2e^{-t} = 1 - 4^{-i}$ for each entry in C_i . For each cluster, even though there are up to 2^i kernel entries, by the triangle inequality it is easy to see that we only need to apply a union bound over a most carefully selected entries (excluding the entries on the diagonal) to meet Eq. (12) with probability at least $1 - 2^i(4^{-i}) = 1 - 2^{-i}$.

Subsequently, applying a union bound over all clusters gives us that with probability at least $1 - \sum_{i=1}^b 2^{-i} = 1 - \frac{1}{2}$, all kernel entries within the t -th cluster satisfy Eq.(12) simultaneously for $1 \leq i \leq b$. \square

Lemma 4. For all $(u; v) \in C^0$, $f(u; v) \leq c(x_u) \oplus c(x_v)$, we have $K(x_u; x_v) < 1 - \frac{1}{2^a}^{\frac{1}{4}}$

where $a = \frac{1}{\log 2} \log \frac{n^4}{n^4 - 4}$ as defined in Lemma 3 above.

Proof. For any $(u; v)$ for which $c(x_u) = i$ and $c(x_v) = j$ and $i \neq j$, we have:

$$\begin{aligned} \frac{1}{2} \log \frac{2^{2^i} (x_u - x_v)^2}{2^{2^i} (c_i - c_j)^2} &= \frac{1}{2} \log \frac{2^{2^i} (x_u - c_i)^2}{2^{2^i} (c_i - c_j)^2} + \frac{1}{2} \log \frac{2^{2^j} (x_v - c_j)^2}{2^{2^j} (c_i - c_j)^2} \\ &= \frac{1}{2} \log \frac{2^{2^i}}{2^{2^i} - 1} + \frac{1}{2} \log \frac{2^{2^j}}{2^{2^j} - 1} \\ &= \frac{1}{2} \log \frac{2^a}{2^a - 1} \\ K(x_u; x_v) &= \exp \left(\frac{1}{2} \log \frac{2^{2^i} (x_u - x_v)^2}{2^{2^i} (c_i - c_j)^2} + \frac{1}{2} \log \frac{2^{2^j} (x_v - c_j)^2}{2^{2^j} (c_i - c_j)^2} \right) \\ &= \exp \left(\frac{1}{2} \log \frac{2^a}{2^a - 1} + \frac{1}{2} \log \frac{2^a}{2^a - 1} \right) < 1 - \frac{1}{2^a}^{\frac{1}{4}} \end{aligned}$$

since for all $(i; j)$, by Condition 3:

$$\frac{1}{2} \log \frac{2^a}{2^a - 1} > \frac{3}{2} \log \frac{2^a}{2^a - 1} \quad (17)$$

This completes our proof for the stated result of Lemma 4. \square

Corollary 1. With probability at least $1 - \frac{1}{2^{(a+b)}}$, there are exactly n entries that are greater than $1 - \frac{1}{2^{(a+b)}}$ where $a = \frac{1}{\log 2} \log \frac{n^4}{n^4 - 4}$. These are the diagonal entries $K(x_u; x_u)$ with $1 \leq u \leq n$.

Proof. Lemma 3 asserts that with probability $1 - \frac{1}{2^{(a+b)}}$, all kernel entries $K(x_u; x_v)$, where $c(x_u) = i$ and $c(x_v) = j$, belong to their respective band $[1 - 2^{-(a+i)}, 1 - 2^{-(a+i+1)}]$ and $[1 - 2^{-(a+j)}, 1 - 2^{-(a+j+1)}]$. When this happens, all in-cluster entries (except the diagonal entries) will have values between $1 - 2^{-a}$ and $1 - 2^{-(a+b)}$ (since there are b bands) and as such, off-cluster entries will either be smaller than $1 - 2^{-a}$ or larger than $1 - 2^{-(a+b)}$. But then Lemma 4 further guarantees that all off-cluster entries are smaller than $1 - 2^{-a}$, following Condition 3. Thus, it follows that the only entries that are larger than $1 - 2^{-(a+b)}$ are the diagonal items and there are exactly n of them. \square

Lemma 5. Let $\mathcal{G}_i = \{(u; v) \mid j = 1 - 2^{-(a+i)}, K^4(x_u; x_v) < 1 - 2^{-(a+i)}\}$. It follows that for each $i \in [1 :: b]$, with probability at least $1 - \frac{1}{2^b} = b$:

$$\sum_{(u;v) \in \mathcal{G}_i} K^2(x_u; x_v) \leq \frac{2}{b} \quad (18)$$

if the kernel approximation $K^0(x_u; x_v)$, $\frac{1}{p} \sum_{t=1}^p K_t(x_u; x_v)$ is formed using at least $\frac{b}{2} \frac{ij}{2^{a+i}} \log \frac{b}{ij} = O\left(\frac{\log^2 n}{2} \log \frac{\log n}{2}\right)$ samples.

Proof. For all $(u; v)$, we have $K_t(x_u; x_v) = \cos(\sum_{i=1}^{t-1} \alpha_i (x_u - x_v))$ where $\alpha_i \sim N(0; 1)$ and,

$$K_t(x_u; x_v) = \cos \sum_{i=1}^{t-1} \alpha_i \frac{x_u - x_v}{\sigma}, \quad \cos z_{uv}^t : \quad (19)$$

Since $\alpha_i \sim N(0; 1)$, z_{uv}^t is then a weighted sum of Gaussian random variables α_i and $N(0; \sum_{i=1}^{t-1} \alpha_i^2)$, where $\sum_{i=1}^{t-1} \alpha_i^2 \sim \chi^2_{t-1}(x_u - x_v)^2$, which in turn implies:

$$\begin{aligned} E[\cos(z_{uv}^t)] &= \exp(-0.5 \sum_{i=1}^{t-1} \alpha_i^2) = K(x_u; x_v); \\ V[\cos(z_{uv}^t)] &= \frac{1}{2} (1 - E[\cos(z_{uv}^t)]^2) = \frac{1}{2} (1 - K^2(x_u; x_v)) \leq \frac{1}{2^{2a+i}}; \end{aligned} \quad (20)$$

where the last inequality follows from the choice of α_i and the definition of the α_i above. Next, applying the Chernoff-Hoeffding inequality and union bounding over $(u; v) \in C$ we have:

$$\begin{aligned} \Pr_{(u;v) \in C} \left[\sum_{i=1}^p \left| \cos(z_{uv}^i) - K(x_u; x_v) \right| \geq \frac{1}{2} \right] &\leq \sum_{i=1}^p \Pr_{(u;v) \in C} \left[\sum_{(u;v) \in C} \left| \cos(z_{uv}^i) - K(x_u; x_v) \right| \geq \frac{1}{4} \right] \\ &\leq \sum_{i=1}^p \sum_{(u;v) \in C} \Pr_{(u;v) \in C} \left[\sum_{(u;v) \in C} \left| \cos(z_{uv}^i) - K(x_u; x_v) \right| \geq \frac{1}{4} \right] \\ &\leq \sum_{i=1}^p \sum_{(u;v) \in C} \exp\left(-\frac{1}{8} \sum_{(u;v) \in C} \frac{1}{2^{2a+i}}\right) \leq \frac{1}{4} \sum_{i=1}^p \sum_{(u;v) \in C} \frac{1}{2^{2a+i}} \end{aligned} \quad (21)$$

Thus, setting $\epsilon = \frac{1}{4} \sum_{i=1}^p \frac{1}{2^{2a+i}}$ and $\delta = \frac{1}{4} \sum_{i=1}^p \frac{1}{2^{2a+i}}$ yields:

$$\Pr_{(u;v) \in C} \left[\sum_{(u;v) \in C} \left| \cos(z_{uv}^i) - K(x_u; x_v) \right| \geq \frac{1}{4} \right] \leq \frac{1}{4} \sum_{i=1}^p \sum_{(u;v) \in C} \frac{1}{2^{2a+i}} \leq \frac{1}{4} \sum_{i=1}^p \frac{1}{2^{2a+i}} \quad (22)$$

where the last inequality follows from the above choice of ϵ . Since $\sum_{i=1}^p \frac{1}{2^{2a+i}} = O\left(\frac{\log^2 n}{2^{2a}}\right)$, we further have $\frac{1}{4} \sum_{i=1}^p \frac{1}{2^{2a+i}} = O\left(\frac{\log^2 n}{2^{2a}}\right)$. \square

Lemma 5 thus establishes a very strong sample complexity $O(\log^2 n \log \log n)$ for approximating all kernel entries within a narrow band of values, which is significantly cheaper than the sample complexity of $O(n^2 \log n)$ we would get if we were to ignore the distribution of kernel values in different bands. This is made clear in Corollary 2 below, which combines Lemmas 3, 4 and 5 to establish an overall sample complexity resulting in only a small approximation loss accumulated over all bands.

Corollary 2. If a kernel approximation K^0 of K is formed such that $K^0(x_u; x_v) = \frac{1}{p} \sum_{i=1}^p K_t(x_u; x_v)$ for all in-cluster entries $(u; v) \in C$ using $p = O(\log^2 n) \log(\log n)$ samples and $K^0(x_{u^0}; x_{v^0}) = 0$ for all off-cluster entries $(u^0; v^0) \in C^0$, then,

$$\|K - K^0\|_F^2 \leq \frac{1}{4} \sum_{(u;v) \in C} \frac{1}{2^{2a+i}};$$

with probability at least $1 - O(\exp(-\log n))$. This immediately guarantees that K^0 is spectrally close to K using the notion of ϵ -closeness (see Definition 1).

Proof. By Lemma 3, with probability $1 - \delta$, $\sum_{i=1}^p \alpha_i = \sum_{i=1}^p \alpha_i$ simultaneously for all $(u; v) \in C$. Thus, applying a union bound over this event and the results obtained in Lemma 5 for all clusters, we have the following bound on the total approximation loss over in-cluster entries $(u; v) \in C$ with probability $1 - 2\delta$:

$$\sum_{(u;v) \in C} \left| \cos(z_{uv}^i) - K(x_u; x_v) \right|^2 \leq \frac{1}{4} \sum_{i=1}^p \sum_{(u;v) \in C} \frac{1}{2^{2a+i}} \quad (23)$$

Furthermore, by Lemma 4, we also have the following bound on the total approximation loss over off-cluster entries $(u^0; v^0) \in C^0$ (which were approximated uniformly by zero):

$$\sum_{(u;v) \in C^0} \left| \cos(z_{uv}^i) - K(x_u; x_v) \right|^2 \leq \frac{3n^2}{4} \sum_{(u;v) \in C^0} \frac{1}{2^{2a+i}} = \frac{3n^2}{4} \sum_{(u;v) \in C^0} \frac{1}{2^{2a+i}} \quad (24)$$

when the last inequality is due to the facts (established in Lemma 4) that $K(x_u; x_v) \leq 1 - \frac{1}{2^a}$ and that $a = \frac{1}{\log 2} \log \frac{n^4}{n^4 - 4}$. Finally, combining these yields:

$$kK - K^{0,2} \leq kK - K^{0,2} = \sum_{(u,v) \in C} \sum_{(u,v) \in C^c} \sum_{(u,v) \in C^c} 2(x_u; x_v) + \sum_{(u,v) \in C^c} \sum_{(u,v) \in C^c} 2(x_u; x_v) - \frac{1}{4} \sum_{(u,v) \in C} 2 + \frac{3}{4} \sum_{(u,v) \in C^c} 2 = \sum_{(u,v) \in C^c} 2 \quad (25)$$

□

B Intermediate Results for Theorem 3

Lemma 6. Let K and K^0 be positive semidefinite matrices in $\mathbb{R}^{n \times n}$ such that $I - K - K^0 \succeq 0$, $Q \succeq K + \frac{1}{2}I$ and $Q^0 \succeq K^0 + \frac{1}{2}I$ for some $\beta > 0$, then:

$$kQ^{0,1} k_2 = (1 - \frac{\beta}{2}) kQ^{1,1} k_2 \quad (26)$$

Proof. By definition of the spectral norm, we have $\exists x \in \mathbb{R}^n$:

$$K - K^0 \succeq -I; \quad (27)$$

which implies

$$\begin{aligned} Q &\succeq K^0 + (\frac{\beta}{2} + \beta)I \\ &\succeq (\frac{\beta}{2} + \beta)K^0 + (\frac{\beta}{2} + \beta)I \\ &= (1 + \frac{\beta}{2}) Q^0; \end{aligned} \quad (28)$$

where β and β denote the Loewner inequality operators. Likewise, by symmetry, we also have:

$$Q^0 \succeq (1 + \frac{\beta}{2}) Q; \quad (29)$$

Let $A = (1 + \frac{\beta}{2})Q^0$ and $B = Q$. Since A and B are symmetric and positive semidefinite, there exist $U; V$ with orthogonal rows and columns and diagonal matrices Λ^0 for which $A = U \Lambda^0 U^T$ and $B = V \Lambda V^T$. We further let $A^{1/2} = U \Lambda^{1/2}$ and $B^{1/2} = V \Lambda^{1/2}$.

Then, we can rewrite Eq. (28) as:

$$\begin{aligned} &A - B \succeq 0 \\ &A^{1/2} B^{1/2} (A - B) B^{1/2} A^{1/2} \succeq 0 \\ &A^{1/2} B^{1/2} A B^{1/2} A^{1/2} \succeq I \\ &A^{1/2} B^{1/2} (B^{1/2} A B^{1/2}) B^{1/2} A^{1/2} \succeq A^{1/2} B^{1/2} B^{1/2} A^{1/2} \\ &A^{1/2} B^{1/2} A^{1/2} B^{1/2} A^{1/2} \succeq I \\ &A^{1/2} (A^{1/2} B^{1/2} A^{1/2}) A^{1/2} \succeq A^{1/2} A^{1/2} \\ &A^{1/2} B^{1/2} A^{1/2} \succeq A^{1/2} \\ &A^{1/2} Q^{1/2} \succeq \frac{2}{2+\beta} Q^{0,1} \\ &(1 + \frac{\beta}{2}) Q^{1/2} \succeq Q^{0,1}; \end{aligned} \quad (30)$$

Again, by symmetry, we can rewrite Eq. 29 as:

$$\begin{aligned} Q^{0,1} &\succeq \frac{2}{2+\beta} Q^{1/2} \succeq (1 - \frac{\beta}{2}) Q^{1/2} \\ &\succeq (1 - \frac{\beta}{2}) Q^{1/2}; \end{aligned} \quad (31)$$

Therefore, we have $kQ^{0,1} k_2 = (1 - \frac{\beta}{2}) kQ^{1,1} k_2$. □

Let $g(\mathbf{x})$ and $g^\theta(\mathbf{x})$ respectively denote the predictive distributions of full GP and the approximated GP pertaining to an arbitrary test input \mathbf{x} . We then state the following lemmas:

Lemma 7. Let \mathbf{K}^θ denote an approximation that is ϵ -close to the original kernel \mathbf{K} . The induced predictive mean of \mathbf{K}^θ is bounded by a factor of $1 \pm \frac{\epsilon}{2}$ times the original predictive mean.

$$\mathbb{E}[g(\mathbf{x})] = \left(1 \pm \frac{\epsilon}{2}\right) \mathbb{E}[g^\theta(\mathbf{x})]; \quad (32)$$

Proof. Let $\mathbf{k} = [k(\mathbf{x}; \mathbf{x}_i)]_{i=1}^n$ where \mathbf{x}_i denotes the i -th training data point. We have:

$$\begin{aligned} \mathbb{E}[g(\mathbf{x})] &= \frac{1}{2} (\mathbf{k} + \mathbf{y})^\top \mathbf{Q}^{-1} (\mathbf{k} + \mathbf{y}) - \mathbf{k}^\top \mathbf{Q}^{-1} \mathbf{k} - \mathbf{y}^\top \mathbf{Q}^{-1} \mathbf{y} \\ &= \frac{1}{2} \left(1 \pm \frac{\epsilon}{2}\right) (\mathbf{k} + \mathbf{y})^\top \mathbf{Q}^{\theta^{-1}} (\mathbf{k} + \mathbf{y}) - \mathbf{k}^\top \mathbf{Q}^{\theta^{-1}} \mathbf{k} - \mathbf{y}^\top \mathbf{Q}^{\theta^{-1}} \mathbf{y} \\ &= \left(1 \pm \frac{\epsilon}{2}\right) \mathbb{E}[g^\theta(\mathbf{x})]; \end{aligned} \quad (33)$$

where the first and third equations follow from adding and subtracting the same terms to the expression of $g(\mathbf{x})$ – see Eq. (2) – while the second equation follows from applying Lemma 6 above. \square

Lemma 8. Let \mathbf{K}^θ denote an approximation that is ϵ -close to the original kernel \mathbf{K} . The induced predictive variance of \mathbf{K}^θ is bounded by a factor of $1 \pm \frac{\epsilon}{2}$ of the original predictive variance up to a constant bias of $\frac{\epsilon}{2}$,

$$\mathbb{V}[g(\mathbf{x})] = \left(1 \pm \frac{\epsilon}{2}\right) \mathbb{V}[g^\theta(\mathbf{x})] \pm \frac{\epsilon}{2}; \quad (34)$$

Proof. Following the definition of the Gaussian kernel, we assume that the signal of the SE (Squared Exponential) kernel is unitary⁸. As such,

$$\begin{aligned} \mathbb{V}[g(\mathbf{x})] &= \mathbf{1}^\top \mathbf{k}^\top \mathbf{Q}^{-1} \mathbf{k} \\ &= \mathbf{1}^\top \left(1 \pm \frac{\epsilon}{2}\right) \mathbf{k}^\top \mathbf{Q}^{\theta^{-1}} \mathbf{k} \\ &= \left(1 \pm \frac{\epsilon}{2}\right) \mathbf{1}^\top \mathbf{k}^\top \mathbf{Q}^{\theta^{-1}} \mathbf{k} \pm \frac{\epsilon}{2} \\ &= \left(1 \pm \frac{\epsilon}{2}\right) \mathbb{V}[g^\theta(\mathbf{x})] \pm \frac{\epsilon}{2}; \end{aligned} \quad (35)$$

where (again) the above equation follows straightforwardly from applying Lemma 6 and standard algebraic manipulation. Lemma 7 and Lemma 8 thus provide an explicit bound on the difference between the original and approximated predictive distributions. We will now establish another bound on the difference between the original and approximated negative log likelihoods (i.e., the training objectives) in Lemma 9 and Lemma 10 below. \square

Lemma 9. Let \mathbf{K}^θ denote an approximation that is ϵ -close to the original kernel \mathbf{K} . Let $\mathbf{Q} = \mathbf{K} + \frac{\epsilon}{2} \mathbf{I}$ and $\mathbf{Q}^\theta = \mathbf{K}^\theta + \frac{\epsilon}{2} \mathbf{I}$. We have:

$$\log j_{\mathbf{Q}^\theta} = \left(1 \pm \frac{\epsilon}{2}\right) \log j_{\mathbf{Q}}; \quad (\mathbf{K}) \quad (36)$$

where the spectral constant $\kappa(\mathbf{K})$ of \mathbf{K} is defined below:

$$\kappa(\mathbf{K}) = \frac{\max(\log(1 + \frac{\epsilon}{2}); \log(1 - \frac{\epsilon}{2}))}{\min(\log(\min(\mathbf{K}) + \frac{\epsilon}{2}); \log(\max(\mathbf{K}) + \frac{\epsilon}{2}))}; \quad (37)$$

⁸This simplifies the analysis and does not restrict the expressiveness of the kernel since we can either normalize the output or absorb it into the length-scales (i.e., the ℓ_i).

Proof. Let $\lambda_1, \lambda_2, \dots, \lambda_n$ and $\lambda_1^0, \lambda_2^0, \dots, \lambda_n^0$ be the eigenvalues of \mathbf{K} and \mathbf{K}^0 respectively. Applying the Courant-Fischer theorem on the result obtained in Lemma 6, we have:

$$\lambda_i^0 + \frac{1}{2} = \lambda_i + \frac{1}{2} \quad ; \quad (38)$$

This implies:

$$\begin{aligned} \log j\mathbf{Q}^0j &= \sum_{i=1}^n \log(\lambda_i^0 + \frac{1}{2}) = \sum_{i=1}^n \log(\lambda_i + \frac{1}{2}) + \log \prod_{i=1}^n \frac{\lambda_i^0 + \frac{1}{2}}{\lambda_i + \frac{1}{2}} \\ &= \sum_{i=1}^n \log(\lambda_i + \frac{1}{2}) + \sum_{i=1}^n \max \left\{ \log \frac{\lambda_i^0 + \frac{1}{2}}{\lambda_i + \frac{1}{2}} ; \log \frac{\lambda_i + \frac{1}{2}}{\lambda_i^0 + \frac{1}{2}} \right\} \\ &= \sum_{i=1}^n \log(\lambda_i + \frac{1}{2}) + \sum_{i=1}^n \max \left\{ \log \frac{\lambda_i^0 + \frac{1}{2}}{\lambda_i + \frac{1}{2}} ; \log \frac{\lambda_i + \frac{1}{2}}{\lambda_i^0 + \frac{1}{2}} \right\} \\ &= \log j\mathbf{Q}j + \sum_{i=1}^n \max \left\{ \log \frac{\lambda_i^0 + \frac{1}{2}}{\lambda_i + \frac{1}{2}} ; \log \frac{\lambda_i + \frac{1}{2}}{\lambda_i^0 + \frac{1}{2}} \right\} \quad ; \quad (\mathbf{K}) \quad (39) \end{aligned}$$

Similarly, by symmetry, we have:

$$\begin{aligned} \log j\mathbf{Q}^0j &= \sum_{i=1}^n \log(\lambda_i^0 + \frac{1}{2}) = \sum_{i=1}^n \log(\lambda_i + \frac{1}{2}) + \sum_{i=1}^n \max \left\{ \log \frac{\lambda_i^0 + \frac{1}{2}}{\lambda_i + \frac{1}{2}} ; \log \frac{\lambda_i + \frac{1}{2}}{\lambda_i^0 + \frac{1}{2}} \right\} \\ &= \sum_{i=1}^n \log(\lambda_i + \frac{1}{2}) + \sum_{i=1}^n \max \left\{ \log \frac{\lambda_i^0 + \frac{1}{2}}{\lambda_i + \frac{1}{2}} ; \log \frac{\lambda_i + \frac{1}{2}}{\lambda_i^0 + \frac{1}{2}} \right\} \\ &= \log j\mathbf{Q}j + \sum_{i=1}^n \max \left\{ \log \frac{\lambda_i^0 + \frac{1}{2}}{\lambda_i + \frac{1}{2}} ; \log \frac{\lambda_i + \frac{1}{2}}{\lambda_i^0 + \frac{1}{2}} \right\} \quad ; \quad (\mathbf{K}) \quad (40) \end{aligned}$$

Together, Eq. (39) and Eq. (40) imply $\log j\mathbf{Q}^0j = \log j\mathbf{Q}j + \sum_{i=1}^n \max \left\{ \log \frac{\lambda_i^0 + \frac{1}{2}}{\lambda_i + \frac{1}{2}} ; \log \frac{\lambda_i + \frac{1}{2}}{\lambda_i^0 + \frac{1}{2}} \right\} \quad ; \quad (\mathbf{K})$. \square

Lemma 10. Let \mathbf{K}^0 denote an approximation that is ϵ -close to the original kernel \mathbf{K} . With λ_i^0 ; (\mathbf{K}) previously defined in Lemma 9, we have:

$$\bar{\ell}^0(\lambda) = \frac{1}{2} \sum_{i=1}^n \max \left\{ \log \frac{\lambda_i^0 + \frac{1}{2}}{\lambda_i + \frac{1}{2}} ; \log \frac{\lambda_i + \frac{1}{2}}{\lambda_i^0 + \frac{1}{2}} \right\} \quad ; \quad (\mathbf{K}) \quad (41)$$

where $\bar{\ell}(\lambda)$ and $\bar{\ell}^0(\lambda)$ respectively denote the negative log likelihood of the full GP and the approximated GP evaluated at the hyper-parameters $\lambda = \text{diag}[\lambda_1^2; \lambda_2^2; \dots; \lambda_n^2]$ as defined previously.

Proof. We have:

$$\begin{aligned} \bar{\ell}^0(\lambda) &= \frac{1}{2} \log j\mathbf{Q}^0j + \frac{1}{2} \mathbf{y}^T \mathbf{Q}^0 \mathbf{1} \mathbf{y} \\ &= \frac{1}{2} \left(\sum_{i=1}^n \log(\lambda_i^0 + \frac{1}{2}) + \mathbf{y}^T \mathbf{Q}^0 \mathbf{1} \mathbf{y} \right) \\ &= \frac{1}{2} \sum_{i=1}^n \max \left\{ \log \frac{\lambda_i^0 + \frac{1}{2}}{\lambda_i + \frac{1}{2}} ; \log \frac{\lambda_i + \frac{1}{2}}{\lambda_i^0 + \frac{1}{2}} \right\} + \frac{1}{2} \log j\mathbf{Q}j + \mathbf{y}^T \mathbf{Q} \mathbf{1} \mathbf{y} \\ &= \frac{1}{2} \sum_{i=1}^n \max \left\{ \log \frac{\lambda_i^0 + \frac{1}{2}}{\lambda_i + \frac{1}{2}} ; \log \frac{\lambda_i + \frac{1}{2}}{\lambda_i^0 + \frac{1}{2}} \right\} + \frac{1}{2} \log j\mathbf{Q}j + \mathbf{y}^T \mathbf{Q} \mathbf{1} \mathbf{y} \quad ; \quad (\mathbf{K}) \quad (42) \end{aligned}$$

Using the result of Lemma 10 above, we can further analyze how the quality of the optimized parameter $\lambda^0 = \arg \max_{\lambda} \bar{\ell}^0(\lambda)$ of the approximated training objective compares to the true optimizer of the original objective function $\lambda = \arg \max_{\lambda} \bar{\ell}(\lambda)$ in Lemma 11 below.

Lemma 11. Let λ and λ^0 denote the optimal hyper-parameters obtained by respectively minimizing the negative log likelihood of the full GP and the approximated GP. We have:

$$\bar{\ell}^0(\lambda^0) = \frac{1}{2} \sum_{i=1}^n \max \left\{ \log \frac{\lambda_i^0 + \frac{1}{2}}{\lambda_i + \frac{1}{2}} ; \log \frac{\lambda_i + \frac{1}{2}}{\lambda_i^0 + \frac{1}{2}} \right\} \quad ; \quad (\mathbf{K}) \quad (43)$$

Proof. By Lemma 10, we have:

$$\begin{aligned} \bar{\ell}^0(\lambda^0) &= \frac{1}{2} \sum_{i=1}^n \max \left\{ \log \frac{\lambda_i^0 + \frac{1}{2}}{\lambda_i + \frac{1}{2}} ; \log \frac{\lambda_i + \frac{1}{2}}{\lambda_i^0 + \frac{1}{2}} \right\} \\ &= \frac{1}{2} \sum_{i=1}^n \max \left\{ \log \frac{\lambda_i^0 + \frac{1}{2}}{\lambda_i + \frac{1}{2}} ; \log \frac{\lambda_i + \frac{1}{2}}{\lambda_i^0 + \frac{1}{2}} \right\} \quad ; \quad (\mathbf{K}) \quad (44) \end{aligned}$$

and

$$\begin{aligned} \psi^0(\theta) &= \frac{1}{2} \max_{\mathbf{K}} \left\{ \log \left(\frac{\mathbf{K}}{\mathbf{K}^0} \right) \right\} \\ &= \frac{1}{2} \max_{\mathbf{K}} \left\{ \log \left(\frac{\mathbf{K}}{\mathbf{K}^0} \right) \right\}. \end{aligned} \quad (45)$$

Together, these results imply $\psi^0(\theta) = \frac{1}{2} \max_{\mathbf{K}} \left\{ \log \left(\frac{\mathbf{K}}{\mathbf{K}^0} \right) \right\}$. \square

Lemma 12. Let $\theta \in (0; 1)$ and let \mathbf{K}^0 denote an approximation of \mathbf{K} for which $\|\mathbf{K} - \mathbf{K}^0\|_2 \leq \epsilon$ with probability at least $1 - \delta$ uniformly over the entire parameter space. Let θ^* and θ^0 denote the optimal hyper-parameters obtained by respectively minimizing the negative log likelihood of the full GP and the approximated GP. Then, with probability $1 - \delta$, the following holds:

$$E[g^0(\mathbf{x}; \theta^0)] = \frac{1}{2} \left(\log \left(\frac{\mathbf{K}^0}{\mathbf{K}} \right) \right) + E[g(\mathbf{x}; \theta^*)] + \lambda(\theta; \theta^0) \quad (46)$$

where $\lambda(\theta; \theta^0)$ and $\lambda(\theta; \theta^0)$ are constant with respect to $\theta; \theta^0$.

Proof. We have:

$$\begin{aligned} E[g(\mathbf{x}; \theta)] &= \mathbf{k}^T \mathbf{Q}^{-1} \mathbf{y} \\ &= \frac{1}{2} (\mathbf{k} + \mathbf{y})^T \mathbf{Q}^{-1} (\mathbf{k} + \mathbf{y}) - \frac{1}{2} \mathbf{k}^T \mathbf{Q}^{-1} \mathbf{k} + \log | \mathbf{Q} | \\ &= \frac{1}{2} \left(\log \left(\frac{\mathbf{K}}{\mathbf{K}^0} \right) \right) + 1 + \sum_{i=1}^n \log \left(\frac{\sigma_i^2}{\sigma_i^0 + \epsilon^2} \right) \end{aligned} \quad (47)$$

On the other hand, we have:

$$\begin{aligned} E[g(\mathbf{x}; \theta^*)] &= \mathbf{k}^T \mathbf{Q}^{-1} \mathbf{y} \\ &= \frac{1}{2} \mathbf{k}^T \mathbf{Q}^{-1} \mathbf{k} + \mathbf{y}^T \mathbf{Q}^{-1} \mathbf{y} \\ &= \frac{1}{2} \left(\log \left(\frac{\mathbf{K}^0}{\mathbf{K}} \right) \right) + 1 + \sum_{i=1}^n \log \left(\frac{\sigma_i^0}{\sigma_i^0 + \epsilon^2} \right) \end{aligned} \quad (48)$$

Thus, we have:

$$E[g(\mathbf{x}; \theta)] = \frac{1}{2} \left(\log \left(\frac{\mathbf{K}}{\mathbf{K}^0} \right) \right) + 1 + \sum_{i=1}^n \log \left(\frac{\sigma_i^2}{\sigma_i^0 + \epsilon^2} \right) \quad (49)$$

and by symmetry:

$$\begin{aligned} E[g^0(\mathbf{x}; \theta^0)] &= \frac{1}{2} \left(\log \left(\frac{\mathbf{K}^0}{\mathbf{K}} \right) \right) + 1 + \sum_{i=1}^n \log \left(\frac{\sigma_i^0}{\sigma_i^0 + \epsilon^2} \right) \\ &= \frac{1}{2} \left(\log \left(\frac{\mathbf{K}^0}{\mathbf{K}} \right) \right) + E[g(\mathbf{x}; \theta^*)] + \lambda(\theta; \theta^0) \end{aligned} \quad (50)$$

where $\lambda(\theta; \theta^0)$ is a constant as defined below:

$$\begin{aligned} \lambda(\theta; \theta^0) &= \max_{\mathbf{K}} \left\{ \log \left(\frac{\mathbf{K}}{\mathbf{K}^0} \right) \right\} \\ \lambda(\theta; \theta^0) &= \sum_{i=1}^n \log \left(\frac{\sigma_i^2}{\sigma_i^0 + \epsilon^2} \right) - \sum_{i=1}^n \log \left(\frac{\sigma_i^0}{\sigma_i^0 + \epsilon^2} \right) \end{aligned} \quad \square$$

C Model Parameterization and Practical Implementation

Our embedding algorithm is based on a VAE implementation where the latent prior, posterior and likelihood of the data generation process are represented via separate mixtures of k Gaussian distributions over a 4-dimensional space. For the latent prior, we set (and fixed) the means of each Gaussian component (i.e., the prior cluster means) at k equidistant points on a 4-dimensional sphere centered at zero with an optimizable radius. For the latent posterior and likelihood, the mean and covariance entries of each component in the mixture are parameterized as outputs of their respective neural networks, which we refer to as Gaussian nets.

In turn, the Gaussian nets are parameterized separately. Each starts with a linear layer comprising of 10 neurons whose outputs are fed simultaneously to two separate hidden (linear) layers with 10 hidden neurons each. Their outputs are then used to form the mean and covariance entries of the corresponding Gaussian component. All neurons are activated by a ReLU unit, and in addition, the (batch) outputs of the first linear layer are also standardized via a learnable 1D batch-norm layer to ensure the stability of batch optimization. The mixing weights that combine such Gaussian nets in the mixtures are also parameterized as the outputs of a linear layer with $k = 8$ neurons where $k = 8$ is also the number of components in our mixture.

The above parameterized latent prior, posterior, and likelihood are then connected in the variational lower-bound (ELBO) as expressed in the first two terms of Eq. (9). This ELBO objective is then combined with two regularization terms weighted with (manually tuned) parameters $\lambda = 8.0$ and $\lambda = 1.2$ as detailed in Eq. (10). The entire function is optimized via gradient descent using the standard Adam optimizer with the default setting implemented in PyTorch.

Once learned, the outputs of the latent posterior were used as the encoded data which were fed as input to our revisited SSGP. For a practical implementation, we also found that additionally passing the encoded data to the latent likelihood generates a reconfigured version of the original data which helps to marginally improve the performance. All of our reported results below are generated with respect to this version of reconfiguration. All of our implementations of GP, SSGP and revisited SSGP that makes use of the output of this reconfiguration process, are also in PyTorch. Our experimental code is released at https://github.com/hqminh/gp_sketch_nips.

D Additional Empirical Results and Visualizations

This section provides additional empirical results and visualizations that complement and corroborate the reported results in the main text. In particular, we provide: (a) a more refined and comprehensive visualization of how our embedding algorithm (Section 3.3) re-configures data across different settings; and (b) an extended comparison with SSGP at different levels of sample complexity when evaluated on middle (10K data points) and large (500K data points) data sets. All data samples used in this section were extracted from the GAS SENSOR dataset [4]⁹.

D.1 The Effect of Data Re-configuration: A Visual Demonstration

This section describes an ablation study to demonstrate the effectiveness of our data re-configuration component (i.e., to approximately meet the practical Conditions 1-3 of our refined analysis). Specifically, we demonstrate this by contrasting the scatter plots of data embeddings (see Fig. 3) before and after reconfiguration using our algorithm in Section 3.3 below. The visualizations are shown for 3 different samples of data, each of which has 10K data points.

For each data sample, its embedding was clustered and re-clustered before and after its reconfiguration. Both clustering processes were generated independently using K-Means to provide an objective visual measurement of the reconfiguration effects of our algorithm.

Observing the above visual excerpts, it appears that after reconfiguration, the clusters across different data samples all became significantly more disengtangled with a visibly increased distance between their cluster centers. This provides conclusive evidence to the data disengtangling effect of our

⁹The entire GAS SENSOR dataset contains approximately 4M data points. However, in the body of this paper we only used a sample of 500K points to conduct our experiments.

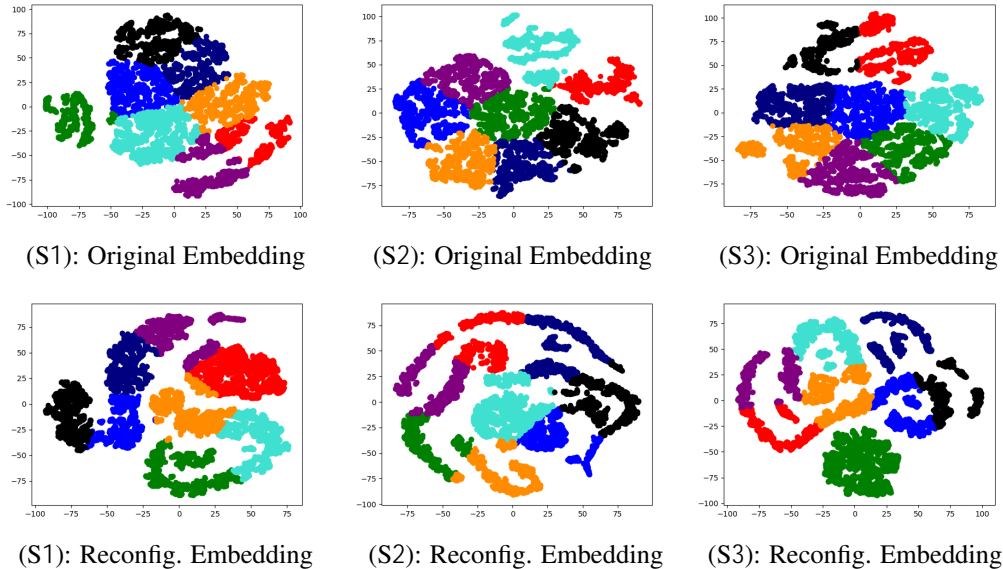


Figure 3: Visualizations of original (top) and reconfigured (bottom) data embeddings for 3 different (randomly selected) data samples annotated with S1, S2 and S3, respectively. Each visual excerpt is annotated with different colors corresponding to the different clusters that the data belong to. All visualizations are generated using T-SNE [41].

embedding algorithm. More importantly, this demonstration further reveals a practical aspect of data that has not been investigated before in the existing literature of GP:

Data (especially experimental data) is often the manifestation of how latent concepts that underlie them were observed and depending on specific parameters of the observation process, these concepts might manifest differently in either more or less useful forms for learning. This raises the question of whether one can reorient the observation process to increase the utility of such data.

In this vein of thought, to address the above question, our data reconfiguration algorithm can be considered to be one potential solution which uses a parameterized construction of a latent space to provide a handle on how to reorient the latent concepts that underlie our data. For an intuitive example, imagine how we would look at the outside world via a narrowed pigeonhole. With different viewing angles, we would perceive the same scene outside differently and apparently, some angles provide a much better perception of that scene (thus, allowing us to interpret the scene more accurately).

In technical terms, such a reorientation is implemented in our algorithm via the regularization of the mixture composition of the latent prior while constraining the entire embedding process to have it reflected on the latent posterior – see Eq. (10) – which was used to encode data into a latent space that exhibits the desired separation effect. Such separation/disentanglement is then shown (empirically) to be richer in information and can be leveraged to improve the sample complexity of SSGP (see Section D.2), thus supporting our theoretical analysis in Appendix A.

D.2 Comparison with SSGP on Large Data

To demonstrate the effectiveness of the data disentanglement in reducing the sample complexity of SSGP, we compare the performance of SSGP and our revisited SSGP (which was instead applied on the reconfigured space of data) at different levels of sample complexity. All results were generated for two different data samples extracted from GAS-SENSOR [4]. One of these (containing 500K data points) is in fact on the same scale of the most extensive datasets used in the GP literature. All performance plots were visualized in Fig. 4 below. For each experiment, the data sample is divided into a train/test partition with an 8-2 ratio. All results were averaged over 5 independent runs.

We see that our revised SSGP consistently achieves better performance than its SSGP counterpart at all complexity levels. In particular, in all cases of the 10K setting, the performance of our revised

