
Appendix: Rates of Convergence for Sparse Variational Gaussian Process Regression

A. Proof Of Lemma 1

Titsias [2014] proves the tighter upper bound,

$$\mathcal{L} \leq \mathcal{L}'_{\text{upper}} := -\frac{N}{2} \log(2\pi) - \frac{1}{2} \log(|\mathbf{Q}_n|) - \frac{1}{2} \mathbf{y}^\top \left(\mathbf{Q}_n + \tilde{\lambda}_{\max} \mathbf{I} \right)^{-1} \mathbf{y}.$$

Subtracting,

$$\begin{aligned} \mathcal{L}'_{\text{upper}} - \mathcal{L}_{\text{lower}} &= \frac{t}{2\sigma_n^2} + \frac{1}{2} \left(\mathbf{y}^\top \left(\mathbf{Q}_n^{-1} - \left(\mathbf{Q}_n + \tilde{\lambda}_{\max} \mathbf{I} \right)^{-1} \right) \mathbf{y} \right). \quad (1) \end{aligned}$$

Since \mathbf{Q}_{ff} is symmetric positive semidefinite, \mathbf{Q}_n is positive definite with eigenvalues bounded below by σ_n^2 . Write, $\mathbf{Q}_n = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^\top$, where \mathbf{U} is unitary and $\mathbf{\Lambda}$ is a diagonal matrix with non-increasing diagonal entries $\gamma_1 \geq \gamma_2 \geq \dots \geq \gamma_N \geq \sigma_n^2$.

We can rewrite the second term (ignoring the factor of one half) in Equation 1 as,

$$\left(\mathbf{U}^\top \mathbf{y} \right)^\top \left(\mathbf{\Lambda}^{-1} - \left(\mathbf{\Lambda} + \tilde{\lambda}_{\max} \mathbf{I} \right)^{-1} \right) \left(\mathbf{U}^\top \mathbf{y} \right).$$

Define, $\mathbf{z} = \left(\mathbf{U}^\top \mathbf{y} \right)$. Since \mathbf{U} is unitary, $\|\mathbf{z}\| = \|\mathbf{y}\|$.

$$\begin{aligned} & \left(\mathbf{U}^\top \mathbf{y} \right)^\top \left(\mathbf{\Lambda}^{-1} - \left(\mathbf{\Lambda} + t\mathbf{I} \right)^{-1} \right) \left(\mathbf{U}^\top \mathbf{y} \right) \\ &= \mathbf{z}^\top \left(\mathbf{\Lambda}^{-1} - \left(\mathbf{\Lambda} + \tilde{\lambda}_{\max} \mathbf{I} \right)^{-1} \right) \mathbf{z} \\ &= \sum_i z_i^2 \frac{\tilde{\lambda}_{\max}}{\gamma_i^2 + \gamma_i \tilde{\lambda}_{\max}} \\ &\leq \|\mathbf{y}\|^2 \frac{\tilde{\lambda}_{\max}}{\gamma_N^2 + \gamma_N \tilde{\lambda}_{\max}}. \end{aligned}$$

The last inequality comes from noting that the fraction in the sum attains a maximum when γ_i is minimized. Since σ_n^2 is a lower bound on the smallest eigenvalue of \mathbf{Q}_n , we have,

$$\mathbf{y}^\top \left(\mathbf{Q}_n^{-1} - \left(\mathbf{Q}_n + \tilde{\lambda}_{\max} \mathbf{I} \right)^{-1} \right) \mathbf{y} \leq \frac{\tilde{\lambda}_{\max} \|\mathbf{y}\|^2}{\sigma_n^4 + \sigma_n^2 \tilde{\lambda}_{\max}}.$$

Lemma 1 follows.

B. KL Divergence Gaussian Distributions

B.1. KL divergence between multivariate Gaussian distributions

We make use of the formula for KL divergences between multivariate Gaussian distributions in our proof of Lemma 2, and the univariate case in Proposition 1.

Recall the KL divergence from $p_1 \sim \mathcal{N}(\mathbf{m}_1, \mathbf{S}_1)$ to $p_2 \sim \mathcal{N}(\mathbf{m}_2, \mathbf{S}_2)$ both of dimension N is given by

$$\begin{aligned} \text{KL}(p_1 \| p_2) &= \frac{1}{2} \left(\text{Tr}(\mathbf{S}_2^{-1} \mathbf{S}_1) + \log \left(\frac{|\mathbf{S}_2|}{|\mathbf{S}_1|} \right) \right. \\ &\quad \left. + (\mathbf{m}_1 - \mathbf{m}_2)^\top \mathbf{S}_2^{-1} (\mathbf{m}_1 - \mathbf{m}_2) - N \right) \geq 0. \quad (2) \end{aligned}$$

The inequality is a special case of Jensen's inequality.

B.2. Proof of Upper Bound in Lemma 2

In the main text we showed,

$$\begin{aligned} \mathbb{E}_y \left[\text{KL}(Q \| \hat{P}) \right] &= \frac{t}{2\sigma_n^2} + \int \mathcal{N}(\mathbf{y}; 0, \mathbf{K}_n) \\ &\quad \times \log \left(\frac{\mathcal{N}(\mathbf{y}; 0, \mathbf{K}_n)}{\mathcal{N}(\mathbf{y}; 0, \mathbf{Q}_n)} \right) d\mathbf{y} \end{aligned}$$

In order to complete the proof, we need to show that the second term on the right hand side is bounded above by $t/(2\sigma_n^2)$. Using Equation 2:

$$\begin{aligned} \mathbb{E}_y \left[\text{KL}(Q \| \hat{P}) \right] &= \frac{t}{2\sigma_n^2} - \frac{N}{2} + \frac{1}{2} \log \left(\frac{|\mathbf{Q}_n|}{|\mathbf{K}_n|} \right) \\ &\quad + \frac{1}{2} \text{Tr}(\mathbf{Q}_n^{-1} \mathbf{K}_n) \\ &\leq \frac{t}{2\sigma_n^2} - \frac{N}{2} + \frac{1}{2} \text{Tr}(\mathbf{Q}_n^{-1} (\mathbf{Q}_n + \tilde{\mathbf{K}}_{\text{ff}})). \quad (3) \end{aligned}$$

The inequality follows from noting the log determinant term is negative, as $\mathbf{K}_n \succ \mathbf{Q}_n$ (i.e. $\mathbf{K}_n - \mathbf{Q}_n$ is positive definite). Simplifying the last term,

$$\begin{aligned} \frac{1}{2} \text{Tr}(\mathbf{I}) + \frac{1}{2} \text{Tr}(\mathbf{Q}_n^{-1} \tilde{\mathbf{K}}_{\text{ff}}) &\leq N/2 + t\lambda_1(\mathbf{Q}_n^{-1})/2 \\ &\leq N/2 + t/(2\sigma_n^2). \end{aligned}$$

The first inequality uses that for positive semi-definite symmetric matrices $\text{Tr}(AB) \leq \text{Tr}(A)\lambda_1(B)$ which is a special case of Hölder's inequality for Schatten norms. The final line uses that the largest eigenvalue of \mathbf{Q}_n^{-1} is bounded above by σ_n^{-2} . Using this in Equation 3 finishes the proof.

B.3. Proof of Proposition 1

Defining $\epsilon = 2KL(q||p)$,

$$\begin{aligned} \epsilon &= \frac{\sigma_1^2 + (\mu_1 - \mu_2)^2}{\sigma_2^2} - \log\left(\frac{\sigma_1^2}{\sigma_2^2}\right) - 1 \\ &\geq x - \log(x) - 1 \end{aligned} \quad (4)$$

where we have defined $x = \frac{\sigma_1^2}{\sigma_2^2}$.

Applying the lower bound $x - \log(x) - 1 \geq (x - 1)^2/2 - (x - 1)^3/3$,

$$\epsilon \geq (x - 1)^2/2 - (x - 1)^3/3.$$

A bound on $|x - 1|$ that holds for all ϵ can then be found with the cubic formula. Under the assumption that $\epsilon < \frac{1}{5}$, we have $x - \log(x) < 1.2$ which implies $x \in [0.493, 1.77]$. For x in this range, we have

$$x - \log(x) - 1 \geq (x - 1)^2/3$$

So,

$$|x - 1| \leq \sqrt{3\epsilon}$$

This proves that,

$$1 - \sqrt{3\epsilon} < \frac{\sigma_1^2}{\sigma_2^2} < 1 + \sqrt{3\epsilon}.$$

From Equation 4 and $x - \log x > 1$,

$$|\mu_1 - \mu_2| \leq \sigma_2 \sqrt{3\epsilon}.$$

Using our bound on the ratio of the variances completes the proof of Proposition 1.

C. Covariances for Interdomain Features

We compute the covariances for eigenvector and eigenfunction inducing features.

C.1. Eigenvector inducing features

Recall we have defined eigenvector inducing features by,

$$u_m = \sum_{i=1}^N w_i^{(m)} f(\mathbf{x}_i).$$

Then,

$$\begin{aligned} \text{cov}(u_m, u_k) &= \mathbb{E} \left[\sum_{i=1}^N w_i^{(m)} f(\mathbf{x}_i) \sum_{j=1}^N w_j^{(k)} f(\mathbf{x}_j) \right] \\ &= \sum_{i=1}^N w_i^{(m)} \sum_{j=1}^N w_j^{(k)} \mathbb{E}[f(\mathbf{x}_i) f(\mathbf{x}_j)] \\ &= \sum_{i=1}^N w_i^{(m)} \sum_{j=1}^N w_j^{(k)} k(\mathbf{x}_i, \mathbf{x}_j). \end{aligned}$$

We now recognize this expression as $\mathbf{w}^{(m)\top} \mathbf{K}_{\mathbf{ff}} \mathbf{w}^{(k)}$. Using the defining property of eigenvectors as well as orthonormality,

$$\text{cov}(u_m, u_k) = \lambda_k (\mathbf{K}_{\mathbf{ff}}) \delta_{m,k}.$$

Similarly,

$$\begin{aligned} \text{cov}(u_m, f(\mathbf{x}_i)) &= \mathbb{E} \left[\sum_{j=1}^N w_j^{(m)} f(\mathbf{x}_j) f(\mathbf{x}_i) \right] \\ &= \sum_{j=1}^N w_j^{(m)} \mathbb{E}[f(\mathbf{x}_j) f(\mathbf{x}_i)] \\ &= \sum_{j=1}^N w_j^{(m)} k(\mathbf{x}_j, \mathbf{x}_i). \end{aligned}$$

This is the i^{th} entry of the matrix vector product $\mathbf{K}_{\mathbf{ff}} \mathbf{w}^{(m)} = \lambda_m (\mathbf{K}_{\mathbf{ff}}) \mathbf{w}_i^{(m)}$.

C.2. Eigenfunction inducing features

Recall we have defined eigenfunction inducing features by,

$$u_m = \int \phi_m(\mathbf{x}) f(\mathbf{x}) p(\mathbf{x}) d\mathbf{x}.$$

Then,

$$\begin{aligned} \text{cov}(u_m, u_k) &= \mathbb{E} \left[\int \phi_m(\mathbf{x}) f(\mathbf{x}) p(\mathbf{x}) d\mathbf{x} \int \phi_k(\mathbf{x}') f(\mathbf{x}') p(\mathbf{x}') d\mathbf{x}' \right] \\ &= \int \phi_m(\mathbf{x}) p(\mathbf{x}) \int \phi_k(\mathbf{x}') \mathbb{E}[f(\mathbf{x}) f(\mathbf{x}')] p(\mathbf{x}') d\mathbf{x}' d\mathbf{x} \\ &= \int \phi_m(\mathbf{x}) p(\mathbf{x}) \int \phi_k(\mathbf{x}') k(\mathbf{x}, \mathbf{x}') p(\mathbf{x}') d\mathbf{x}' d\mathbf{x}. \end{aligned}$$

The expectation and integration may be interchanged by Fubini's theorem, as both integrals converge absolutely since $p(\mathbf{x})$ is a probability density, the $\phi_m(\mathbf{x})$ are in $L^2(\mathcal{X})_p \subset L^1(\mathcal{X})_p$ and k is bounded.

We may then apply the eigenfunction property to the inner integral and orthonormality of eigenfunctions to the result yielding,

$$\text{cov}(u_m, u_k) = \lambda_k \int \phi_k(\mathbf{x}) \phi_m(\mathbf{x}) p(\mathbf{x}) d\mathbf{x} = \lambda_k \delta_{m,k}.$$

With similar considerations,

$$\begin{aligned} \text{cov}(u_m, f(\mathbf{x}_i)) &= \mathbb{E} \left[\int \phi_m(\mathbf{x}) f(\mathbf{x}) f(\mathbf{x}_i) p(\mathbf{x}) d\mathbf{x} \right] \\ &= \int \phi_m(\mathbf{x}) \mathbb{E}[f(\mathbf{x}) f(\mathbf{x}_i)] p(\mathbf{x}) d\mathbf{x} \\ &= \lambda_m \phi_m(\mathbf{x}_i). \end{aligned}$$

Algorithm 1 MCMC algorithm for sampling ϵ k-DPP (A)

Input: Training inputs $\mathbf{X} = \{\mathbf{x}_i\}_{i=1}^N$, number of points to choose, M , kernel k .

Returns: A sample of M inducing points drawn proportional to the determinant of $\mathbf{K}_{Z,Z}$

Initialize M columns greedily in an iterative fashion call this set S_0 .

for $r < R$ **do**

 Sample i uniformly from S_r and j uniformly from $\mathbf{X} \setminus S_r$. Define $T = S \setminus \{i\} \cup \{j\}$,

 Compute $p_{i \rightarrow j} := \frac{1}{2} \min\{1, \det(\mathbf{K}_T)/\det(\mathbf{K}_{S_r})\}$

 With probability $p_{i \rightarrow j}$ $S_{r+1} = T$ otherwise, $S_{r+1} = S$

end for

Return: S_R

D. Discrete k-DPPs

D.1. Proof of Corollary 1 from Lemma 3

$$\begin{aligned} \mathbb{E}_{Z \sim \nu}[t] &= \mathbb{E}_{Z \sim \mu}[t] + (\mathbb{E}_{Z \sim \nu}[t] - \mathbb{E}_{Z \sim \mu}[t]) \\ &\leq (M+1) \sum_{M+1}^N \lambda(\mathbf{K}_{\text{ff}}) + \sum_{Z \in \binom{N}{M}} (\mu(Z) - \nu(Z))t(Z) \\ &\leq (M+1) \sum_{M+1}^N \lambda(\mathbf{K}_{\text{ff}}) + N\nu \sum_{Z \in \binom{N}{M}} |\mu(Z) - \nu(Z)| \\ &\leq (M+1) \sum_{M+1}^N \lambda(\mathbf{K}_{\text{ff}}) + 2N\nu\epsilon. \end{aligned}$$

The first inequality follows from Lemma 3. The second uses the triangle inequality replace $t(Z)$ with a bound on its maximum. The final line uses one of the definitions of total variation distance for discrete random variables.

D.2. Sampling Approximate k-DPPs

Belabbas & Wolfe [2009] proposed using the Metropolis method for approximate sampling from a k-DPP. Several recent works have shown that a natural Metropolis algorithm on k-DPPs mixes quickly. In particular, Anari et al. [2016] considers the following algorithm:

Theorem 1 (Anari et al. [2016], Theorem 2). *Let A denote algorithm 1. Let ν^R denote the distribution induced by R steps of A . Let $R(\epsilon)$ denote the minimum R such that $\|\mu - \nu^R\|_{TV} < \epsilon$, where μ is a k-DPP on some kernel matrix \mathbf{K}_{ff} . Then*

$$R(\epsilon) \leq NM \log \left(\frac{\binom{N}{M} M!}{\epsilon} \right) \leq NM^2 \log N + NM \log \frac{1}{\epsilon}.$$

Taking ϵ to be any fixed inverse power of N , (i.e. $\epsilon = N^{-\gamma}$,

will make the second term $\mathcal{O}(NM \log(N))$, while by taking γ large (e.g. greater than 2), we can make $2N\nu\epsilon$ small.

The total cost of the algorithm is determined by the cost of the greedy initialization, plus $R(\epsilon)$ times the per iteration cost of the algorithm. A naive implementation of the greedy initialization requires $\mathcal{O}(NM^4)$ time and $\mathcal{O}(NM)$ memory, simply by computing the determinant of each of the $N - m$ possible ways to extend the current subset (faster implementations are possible, but this suffices for our purposes). We assume that this is implemented in such a way that at the end of the initialization we have access to $\det(\mathbf{K}_{S_0})$ and a Cholesky factorization $S_0 = \mathbf{L}_0 \mathbf{L}_0^\top$.

We take as an inductive hypothesis that at iteration r of the algorithm, we know $\det(\mathbf{K}_{S_r})$ and a Cholesky factorization, $\mathbf{K}_{S_r} = \mathbf{L}_r \mathbf{L}_r^\top$. We then need to show we can compute a Cholesky factorization and determinant of \mathbf{K}_T in $\mathcal{O}(M^2)$. Given the Cholesky factorization of T , $\det(\mathbf{K}_T)$ can be computed $\mathcal{O}(M)$ as it is the product of the square of the diagonal elements. It therefore remains to consider the calculation of \mathbf{L}_T , a Cholesky factor of \mathbf{K}_T .

The computation of L_T proceeds in two steps: first we compute $\mathbf{L}_{S \setminus i}$ using \mathbf{L}_S .

$$\mathbf{L}_S = \begin{bmatrix} \mathbf{L}_{1,1} & 0 & 0 \\ \ell_{2,1} & \ell_{2,2} & 0 \\ \mathbf{L}_{3,1} & \ell_{3,2} & \mathbf{L}_{3,3} \end{bmatrix}, \quad \mathbf{K}_S = \begin{bmatrix} \mathbf{K}_{1,1} & \mathbf{k}_{2,1} & \mathbf{K}_{1,3} \\ \mathbf{k}_{2,1} & k(\mathbf{x}_i, \mathbf{x}_i) & \mathbf{k}_{2,3} \\ \mathbf{K}_{3,1} & \mathbf{k}_{3,2} & \mathbf{K}_{3,3} \end{bmatrix}$$

A direct calculation shows,

$$L_{S \setminus i} = \begin{bmatrix} \mathbf{L}_{1,1} & 0 \\ \mathbf{L}_{3,1} & \mathbf{L}'_{3,3} \end{bmatrix}$$

where $\mathbf{L}'_{3,3} \mathbf{L}'_{3,3}{}^\top = \mathbf{L}_{3,3} \mathbf{L}_{3,3}{}^\top + \ell_{3,2} \ell_{3,2}{}^\top$. This is a rank one update to a Cholesky factorization, and can be performed using standard methods in $\mathcal{O}(M^2)$.

We now need to extend a Cholesky factorization from $S \setminus i$ to T , which involves adding a row.

$$\mathbf{L}_T = \begin{bmatrix} \mathbf{L}_{S \setminus i} & 0 \\ \mathbf{c} & d \end{bmatrix}$$

with $\mathbf{c} = \mathbf{L}_{S \setminus i}^{-1} \mathbf{k}_{S \setminus i, j}$ and $d = \sqrt{k(\mathbf{x}_j, \mathbf{x}_j) - \mathbf{k}_{S \setminus i, j}^\top \mathbf{L}_{S \setminus i}^{-1} \mathbf{L}_{S \setminus i}^{-1} \mathbf{k}_{S \setminus i, j}}$. All of these calculations can be computed in $\mathcal{O}(M^2)$ completing the proof of the per iteration cost.

E. Proof of Corollaries

E.1. Corollary 2

From Theorem 3, with probability $1 - \delta$ and this choice of ϵ

$$\text{KL}(Q \parallel \hat{P}) \leq \frac{C(M+1)}{2\sigma_n^2 \delta} \left(1 + \frac{\|\mathbf{y}\|_2^2}{\sigma_n^2} \right) + N^{-\gamma} (R/\sigma_n^2 + 1/N) \quad (5)$$

Take $M = \frac{(3+\gamma)\log(N)+\log D}{\log(B^{-1})}$. If $M \geq N$ the KL-divergence is zero and we are done. Otherwise, $C(M+1) < N^2 \sum_{i=M+1}^{\infty} \lambda_i$. By the geometric series formula,

$$\sum_{i=M+1}^{\infty} \lambda_i = v \frac{\sqrt{2a}}{\sqrt{A}} \frac{B^M}{1-B}$$

Now, $B^M = N^{-3-\gamma} D^{-1}$, so

$$\sum_{i=M+1}^{\infty} \lambda_i = \delta N^{-3-\gamma},$$

implying $\frac{C(M+1)}{2\delta\sigma_n^2} < N^{-1-\gamma}$. Using this in eq. 5 completes the proof.

E.2. Corollary 3

It is sufficient to consider the case of isotropic kernels and input distributions.¹ From [Seeger et al., 2008] in the isotropic case (i.e. $B_i = B_j =: B$ for all $i, j \leq D$),

$$\lambda_{s+D-1} \leq \left(\frac{2a}{A}\right)^{D/2} B^{s^{1/D}}.$$

Define $\tilde{M} = M + D - 1$, to be the number of features used for inference.

$$\sum_{s=\tilde{M}+1}^{\infty} \lambda_s \leq \left(\frac{2a}{A}\right)^{D/2} \sum_{s=M+1}^{\infty} B^{s^{1/D}} \quad (6)$$

$$< \left(\frac{2a}{A}\right)^{D/2} \int_{s=M}^{\infty} B^{s^{1/D}} ds =: \mathcal{I}. \quad (7)$$

In the second line, we use that $B < 1$, so $B^{s^{1/D}}$ obtains its minimum on the interval $s \in [s', s'+1]$ at the right endpoint (i.e. monotonicity). We now define $\alpha = -\log(B)$. So,

$$\mathcal{I} = \left(\frac{2a}{A}\right)^{D/2} \int_{s=M}^{\infty} \exp(-\alpha s^{1/D}) ds \quad (8)$$

$$= \left(\frac{2a}{A}\right)^{D/2} \alpha^{-D} D \int_{t=\alpha M^{1/D}}^{\infty} e^{-t} t^{D-1} dt \quad (9)$$

In the second line we made the substitution $t = \alpha s^{1/D}$, so $ds = \alpha^{-D} D t^{D-1}$. We now recognize,

$$\int_{t=\alpha M^{1/D}}^{\infty} e^{-t} t^{D-1} dt$$

¹For the general case, the eigenvalues can be bounded above by constant times the eigenvalues of an operator with an isotropic kernel with all lengthscales equal to the shortest kernel lengthscale and the input density standard deviation set to the largest standard deviation of any one-dimensional marginal of $p(\mathbf{x})$.

as an incomplete gamma function, $\Gamma(D, \alpha M^{1/D})$. From Gradshteyn & Ryzhik [2014, 8.352],

$$\Gamma(D, y) = (D-1)! e^{-y} \sum_{k=0}^{D-1} \frac{y^k}{k!}$$

So,

$$\mathcal{I} = \left(\frac{2a}{A}\right)^{D/2} \alpha^{-D} D! e^{-\alpha M^{1/D}} \sum_{k=0}^{D-1} \frac{\alpha^k M^{k/D}}{k!} \quad (10)$$

As M grows as a function of N and D is fixed, for N large $D \leq \alpha M^{1/D}$. This implies that the largest term in the sum on the right hand side is the final term, so

$$\mathcal{I} \leq \left(\frac{2a}{A}\right)^{D/2} \alpha^{-1} D^2 e^{-\alpha M^{1/D}} M^{(D-1)/D}. \quad (11)$$

Choose $M = \frac{1}{\alpha} \log\left(N^{\gamma'} \left(\frac{2a}{A}\right)^{D/2} D^2 \alpha^{-1}\right)^D$. Then

$$\left(\frac{2a}{A}\right)^{D/2} \alpha^{-1} D^2 e^{-\alpha M^{1/D}} = N^{-\gamma'}$$

and

$$M^{(D-1)/D} < M = \frac{1}{\alpha} \log\left(N^{\gamma'} \left(\frac{2a}{A}\right)^{D/2} D^2 \alpha^{-1}\right)^D,$$

so

$$\mathcal{I} \leq \frac{1}{\alpha} \log\left(N^{\gamma'} \left(\frac{2a}{A}\right)^{D/2} D^2 \alpha^{-1}\right)^D N^{-\gamma'}.$$

For any fixed D , for this choice of M for any $\varepsilon > 0$ for N large,

$$\mathcal{I} = \mathcal{O}\left(N^{-\gamma'+\varepsilon}\right).$$

By choosing $\gamma' > 3 + \varepsilon'$, for some fixed $\varepsilon' > 0$ the proof is complete, using a similar argument as the one used in the proof of the previous corollary.

Note that using the bound proven in Theorem 2 (the tightest of our bounds) the exponential scaling in dimension is unavoidable. If both k and $p(\mathbf{x})$ are isotropic, then the eigenvalue $\left(\frac{2a}{A}\right)^D B^m$ appears $\binom{m+D-1}{D-1}$ times. This follows from noting that this is the number of ways to write m as a sum of D non-negative integers. Using an identity and standard lower bound for binomial coefficients $\sum_{i=1}^K \binom{m+D-1}{D-1} = \binom{K+D}{D} \geq \left(\frac{K+D}{D}\right)^D > C(D)K^D$. for some constant depending on D , $C(D)$. Using Theorem 2 we need to choose M such that $\lambda_M = \mathcal{O}(1/N)$. This means choosing $K \gg \log(N)$ in the sum above, leading to at least $\alpha \log^D(N)$ features being needed for some constant α . The constant in this lower bound decays rapidly with

D , while the constant in the upper bound does not. Better understanding this gap is important for understanding the performance of sparse Gaussian process approximations in high dimensions.

If the data actually lies on a lower dimensional manifold, we conjecture the scaling depends mainly on the dimensionality of the manifold. In particular, if the manifold is linear and axis-aligned, then the kernel matrix only depends on distances along the manifold (not in the space it is embedded in) so the eigenvalues will not be effected by the higher dimensional embedding. We conjecture that similar properties are exhibited when the data manifold is nonlinear.

F. Smoothness and Sacks-Ylvisaker Conditions

In many instances the precise eigenvalues of the covariance operator are not available, but the asymptotic properties are well understood. A notable example is when the data is distributed uniformly on the unit interval. If the kernel satisfies the Sacks-Ylvisaker conditions of order r :

- $k(x, x')$ is r -times continuously differentiable on $[0, 1]^2$. Moreover, $k(x, x')$ has continuous partial derivatives up to order $r+2$ times on $(0, 1)^2 \cap (x > x')$ and $(0, 1)^2 \cap (x < x')$. These partial derivatives can be continuously extended to the closure of both regions.
- Let L denote $k^{(r,r)}(x, x')$, L_+ denote the restriction of L to the upper triangle and L_- the restriction to the lower triangle, then $L_+^{(1,0)} < L_-^{(1,0)}$ on the diagonal $x = x'$.
- $L_+^{(2,0)}(s, \cdot)$ is an element of the RKHS associated to L and has norm bounded independent of s .

Notably, Matérn half integer kernels of order $r + 1/2$ meet the S-Y condition of order r . See [Ritter et al., 1995] for a more detailed explanation of these conditions and extensions to the multivariate case.

References

- Anari, N., Gharan, S. O., and Rezaei, A. Monte Carlo Markov Chain Algorithms for Sampling Strongly Rayleigh Distributions and Determinantal Point Processes. In *29th Annual Conference on Learning Theory (COLT)*, pp. 103–115, 2016.
- Belabbas, M.-A. and Wolfe, P. J. Spectral Methods in Machine Learning and new Strategies for very Large Datasets. In *Proceedings of the National Academy of Sciences (PNAS)*, volume 106, pp. 369–374, 2009.
- Gradshteyn, I. S. and Ryzhik, I. M. *Table of Integrals, Series, and Products*. Academic press, 2014.
- Ritter, K., Wasilkowski, G. W., and Woźniakowski, H. Multivariate Integration and Approximation for Random Fields Satisfying Sacks-Ylvisaker Conditions. In *The Annals of Applied Probability*, volume 5, pp. 518–540. Institute of Mathematical Statistics, 1995.
- Seeger, M. W., Kakade, S. M., and Foster, D. P. Information Consistency of Nonparametric Gaussian Process Methods. In *IEEE Transactions on Information Theory*, volume 54, pp. 2376–2382. IEEE, 2008.
- Titsias, M. K. Variational Inference for Gaussian and Determinantal Point Processes. In *Workshop on Advances in Variational Inference (NIPS)*, December 2014.