

Gaussian Processes and Kernel Methods: A Review on Connections and Equivalences

Motonobu Kanagawa¹, Philipp Hennig¹,
Dino Sejdinovic², and Bharath K Sriperumbudur³

¹University of Tübingen and Max Planck Institute for Intelligent Systems,
Max-Planck-Ring 4, 72076 Tübingen, Germany
{motonobu.kanagawa, ph}@tue.mpg.de

²Department of Statistics, University of Oxford,
24-29 St Giles', Oxford OX1 3LB, UK
dino.sejdinovic@stats.ox.ac.uk

³Department of Statistics, Pennsylvania State University,
University Park, PA 16802, USA
bks18@psu.edu

July 10, 2018

Abstract

This paper is an attempt to bridge the conceptual gaps between researchers working on the two widely used approaches based on positive definite kernels: Bayesian learning or inference using Gaussian processes on the one side, and frequentist kernel methods based on reproducing kernel Hilbert spaces on the other. It is widely known in machine learning that these two formalisms are closely related; for instance, the estimator of kernel ridge regression is identical to the posterior mean of Gaussian process regression. However, they have been studied and developed almost independently by two essentially separate communities, and this makes it difficult to seamlessly transfer results between them. Our aim is to overcome this potential difficulty. To this end, we review several old and new results and concepts from either side, and juxtapose algorithmic quantities from each framework to highlight close similarities. We also provide discussions on subtle philosophical and theoretical differences between the two approaches.

Contents

1	Introduction	2
1.1	Contents of the Paper	4
1.2	Notation and Definitions	6
2	Gaussian Processes and RKHSs: Preliminaries	6
2.1	Positive Definite Kernels	7

2.2	Gaussian Processes	8
2.3	Reproducing Kernel Hilbert Spaces	10
2.4	A Spectral Characterization for RKHSs Associated with Shift-Invariant Kernels	12
3	Connections between Gaussian Process and Kernel Ridge Regression	14
3.1	Gaussian Process Regression and Interpolation	15
3.2	Kernel Ridge Regression and Kernel Interpolation	18
3.3	Equivalences in Regression and Interpolation	21
3.4	Error Estimates: Posterior Variance and Worst-Case Error	23
3.5	A Weight Vector Viewpoint of Regularization and the Additive Noise Assumption	27
4	Hypothesis Spaces: Do Gaussian Process Draws Lie in an RKHS?	28
4.1	Characterizations via Orthonormal Expansions	29
4.1.1	Mercer’s Theorem	29
4.1.2	Mercer Representation of RKHSs	30
4.1.3	Karhunen-Loève Expansion of Gaussian Processes	31
4.2	Sample Path Properties and the Zero-One Law	32
4.3	Powers of RKHSs as GP Sample Spaces	35
4.4	Examples of Sample Path Properties	37
5	Convergence and Posterior Contraction Rates in Regression	38
5.1	Convergence Rates for Gaussian Process and Kernel Ridge Regression	38
5.2	Upper-bounds and Contraction Rates for Posterior Variance	41
6	Integral Transforms	42
6.1	Maximum Mean Discrepancy: Worst Case and Average Case Errors	43
6.2	Sampling and Numerical Integration	46
6.3	Kernel Mean Shrinkage Estimator and Its Bayesian Interpretation	49
6.4	Gaussian Process Interpretation of Hilbert Schmidt Independence Criterion . .	52
7	Conclusions	54
A	Proofs	55
A.1	Proof of Lemma 3.9	55
A.2	Proof of Corollary 4.13	55
A.3	Proof of Proposition 6.4	57

1 Introduction

In machine learning, two nonparametric approaches based on positive definite kernels have been widely used for the purpose of modeling nonlinear functional relationships. On the one side, there is Bayesian machine learning with Gaussian processes (GP), which models a problem at hand probabilistically and produces a posterior distribution for an unknown function of interest. On the other, frequentist kernel methods with reproducing kernel Hilbert spaces (RKHS) usually take a decision theoretic approach by defining a loss function and optimizing the empirical risk. These two approaches have been shown to be practically powerful and

theoretically sound, and have found a wide range of practical applications in dealing with nonlinear phenomena.

It is well known that the two approaches are intimately connected. The most notable example is that, if both use the same kernel, the posterior mean of Gaussian process regression equals the estimator of kernel ridge regression [Kimeldorf and Wahba, 1970]. Another connection is between Bayesian quadrature [O’Hagan, 1991] and kernel herding [Chen et al., 2010], which are in fact equivalent approaches to numerical integration or deterministic sampling [Huszár and Duvenaud, 2012]. These equivalences arise from the more fundamental connection that the notion of positive definite kernels is leveraged both in Gaussian processes as covariance functions, and in RKHSs as reproducing kernels.

There are also less deeply studied connections between the Bayesian and the frequentist approaches. The posterior variance plays a fundamental role in the Bayesian approach, where it quantifies uncertainty over latent quantities of interest. As we show in Section 3.4, posterior variance can be interpreted as a worst case error in an RKHS. This frequentist interpretation is much less widely known, and less well understood. It is rarely mentioned in the literature on frequentist kernel methods, and some of its potential applications therein may have been missed.

The two approaches also subtly differ in how they define hypothesis spaces, which is a core aspect of statistical methods. Consider for instance the regression problem, which involves a hypothesis space for the unknown regression function. The Bayesian approach defines a hypothesis space through a GP prior distribution, treating the true function as a random function. The support of the GP then expresses the knowledge or belief over the truth, and the probability mass expresses the degree of belief. On the other hand, the frequentist approach expresses one’s prior knowledge or belief by assuming the truth belongs to an RKHS or can be approximated well by functions in the RKHS. Even though the use of the same kernel leads to similar structural assumptions about the function of interest in both approaches, e.g., periodicity or smoothness, there is a fundamental modeling difference, because the support of a GP is *not* identical to the corresponding RKHS (e.g., Lukić and Beder [2001, Corollary 7.1]; see also Section 4.2). In fact, sample paths of the GP fall outside the RKHS of the covariance kernel with probability one. This fact might give researchers an impression that differences outweigh the similarities and that the known connections between the Bayesian and frequentist approaches are rather superficial. However, as we show in Sections 4 and 5, a closer look reveals further deep connections.

This text reviews known, and establishes new, equivalences between the Bayesian and frequentist approaches. Our aim is to help researchers working in both fields gain mutual understanding, and be able to seamlessly transfer results in either side to the other. In fact there are some quantities that are almost exclusively studied and utilized on one side of the debate, and this may highlight interesting directions for the other community. Our second motivation is that, while the connections between the two approaches are found and mentioned individually in papers or books, we are not aware of thorough texts focusing specifically on this topic from a modern perspective. We thus aim to collect a short yet systematic overview of the connections. Finally, we also hope that this text offers a short pedagogical introduction to researchers and students who are new to and interested in either of the two fields.

1.1 Contents of the Paper

The principal results mentioned in the later text can be summarized as follows.

Section 2: Gaussian Processes and RKHSs: Preliminaries As a starting point, we review basic definitions and properties of GPs and RKHSs with illustrative examples. We also provide a characterization of RKHSs based on Fourier transforms of kernels, which helps the reader to understand the structure of RKHSs in terms of smoothness of functions.

Section 3: Connections between Gaussian Process and Kernel Ridge Regression Regression is arguably one of the most basic and practically important problems in machine learning and statistics. We consider Gaussian process regression and kernel ridge regression, and discuss equivalences between the two methods. As mentioned above, it is well known that the posterior mean in GP-regression coincides with the estimator of kernel ridge regression. We furthermore show that there is a frequentist interpretation for posterior variance in GP-regression, as a worst case error in kernel ridge regression. In this sense, average-case and worst-case error are equivalent in the least-squares setting, which is key to understanding the connections between the Bayesian and frequentist approaches.

We also discuss the role of additive Gaussian noise in GP-regression and regularization in kernel ridge regression, showing that they are essentially equivalent as a mechanism to make regressors smoother. We then discuss the noise-free setting, where regression becomes interpolation. Here the equivalence between the two approaches can be useful: an upper-bound on posterior variance is derived, transferring a result from the frequentist literature on scattered data approximation to the Bayesian setting, as shown in Section 5.2.

Section 4: Hypothesis Spaces: Do Gaussian Process Draws Lie in an RKHS? We review the properties of GPs and RKHSs as hypothesis spaces, that is, as a way of expressing prior knowledge or belief. We begin with characterizations of GPs and RKHSs by orthonormal expansions, known respectively as the Mercer representation and the Karhunen-Loève expansion. These characterizations allow us to phrase quantities of interest in terms of eigenvalues and eigenfunctions of an integral operator defined by the kernel. We then discuss previous results of Driscoll [1973], Lukić and Beder [2001] providing a necessary and sufficient condition for a given GP to be a member of a given RKHS (which can be different from the RKHS associated with the covariance kernel of the GP). This implies that, while GP sample paths are almost surely outside of the corresponding RKHS, they lie in a function space “slightly larger” than the RKHS, which is itself a certain RKHS [Steinwart and Scovel, 2012, Steinwart, 2017]. In this sense, the Bayesian prior and the frequentist hypothesis space, while not identical, are arguably closer to each other than is often acknowledged.

Section 5: Convergence and Posterior Contraction Rates in Regression We compare convergence properties of GP-regression and kernel ridge regression. Specifically, we show that convergence rates for GP-regression obtained in van der Vaart and van Zanten [2011] can be recovered from those for kernel ridge regression obtained in Steinwart et al. [2009], considering the situation where a regression function is assumed to have a finite degree of smoothness. Since the GP prior is supported on a slightly larger space than the RKHS, to recover convergence rates matching that of GP regression, we need to assume that the

regression function belongs to this slightly larger space. Even in this case, one can obtain a convergence rate for kernel ridge regression, thanks to the approximation capability of the RKHS. That is, the regression function can be approximated well by functions in the RKHS, with the accuracy of approximation determined by the choice of a regularization constant. Interestingly, the asymptotically optimal schedule of regularization constants translates to the assumption in GP regression that noise variance remains constant with increasing sample size. In this sense, a Bayesian assumption of additive noise is related to regularization in the frequentist approach.

Section 6: Integral Transforms This section deals with somewhat more exotic topics than regression, where connections between the Bayesian and frequentist literature have not been studied as deeply. We discuss integral transforms of (probability) measures with kernels, a framework known as *kernel mean embeddings of distributions* [Smola et al., 2007]. This approach provides a nonparametric way of representing probability distributions, and of measuring a distance between them. The former are called *kernel means*, and the latter the *maximum mean discrepancy* (MMD). These have been widely used in machine learning, and interested readers may have a look at the recent survey by Muandet et al. [2017].

While the MMD is characterized as the *worst-case* error of integrals with respect to functions in an RKHS, it can also be characterized as the *average-case error* of integrals with respect to draws from the corresponding GP [Ritter, 2000, Corollary 7 in p.40]. This viewpoint provides an alternative way to understand kernel embeddings in the language of Bayesian quadrature for Bayesians who are familiar with GPs but not with RKHSs, and vice versa. We also review a shrinkage estimator for kernel means proposed by Muandet et al. [2016] and the corresponding GP-based Bayesian interpretation [Flaxman et al., 2016]. We then discuss the problem of sampling or numerical integration, for which GPs and kernel methods have also played fundamental roles in the form of integral transforms [O’Hagan, 1991, Hickernell, 1998, Briol et al., 2018, Dick et al., 2013].

Finally, we study the connections between GPs and kernel methods as applied to the problem of measuring dependence between random variables. We consider the Hilbert-Schmidt independence criterion (HSIC), a kernel-based measure for dependency between two random variables [Gretton et al., 2005], which has a wide range of applications including independence testing, variable selection and causal discovery. HSIC is defined in terms of RKHSs via the cross-covariance operator or via the joint kernel embedding of two random variables; this definition makes HSIC difficult to interpret without close familiarity with RKHSs. We give an alternative formulation of HSIC in terms of GPs, which is to the best of our knowledge novel, and recovers Brownian distance covariance proposed by Székely and Rizzo [2009]. We believe this result provides an intuitive explanation for people whose backgrounds are from GPs about why HSIC is a sensible dependency measure.

Related Literature We collect here a few key related literature on GPs and kernel methods that may be helpful for further reading. Our aim is the closest in spirit to Berlinet and Thomas-Agnan [2004], who collected classic results on the use of RKHSs in statistics and probability; these include the results by Kolmogorov [1941], Parzen [1961], Matheron [1962], Kimeldorf and Wahba [1970], Larkin [1972], who made the earliest contributions to the field. Other related books and monographs include Wahba [1990], Adler [1990], Janson [1997], Stein

[1999], Ritter [2000], Schölkopf and Smola [2002], Wendland [2005], Schaback and Wendland [2006], Rasmussen and Williams [2006], Adler and Taylor [2007], Steinwart and Christmann [2008], van der Vaart and van Zanten [2008], Novak and Wózniaowski [2008, 2010], Stuart [2010], Scheuere et al. [2013], Hennig et al. [2015], Muandet et al. [2017].

1.2 Notation and Definitions

We collect the notation and definitions that will be used throughout the paper.

Basics For a matrix (or a vector) M , its transpose is denoted by M^\top . Let \mathbb{N} be the set of natural numbers, $\mathbb{N}_0 := \mathbb{N} \cup \{0\}$ and \mathbb{N}_0^d be the d -dimensional Cartesian product of \mathbb{N}_0 with $d \in \mathbb{N}$. For a multi-index $\alpha := (\alpha_1, \dots, \alpha)^\top \in \mathbb{N}_0^d$, define $|\alpha| := \sum_{i=1}^d \alpha_i$. \mathbb{R} denotes the real line, \mathbb{R}^d the d -dimensional Euclidean space for $d \in \mathbb{N}$, and $\|\cdot\|$ the Euclidean norm. For $\alpha \in \mathbb{N}_0^d$ and a function f defined on \mathbb{R}^d , let $\partial^\alpha f$ and $D^\alpha f$ be the α -th partial derivative and the α -th partial weak derivative, respectively. For $i, j \in \mathbb{N}$, we define $\delta_{ij} \in \{0, 1\}$ as $\delta_{ij} = 1$ if $i = j$, and $\delta_{ij} = 0$ otherwise.

Probability For a random variable x and a probability distribution P , writing $x \sim P$ means that x has distribution P . $\mathbb{E}[\cdot]$ denotes the expectation of the argument in the bracket, with respect to a random variable concerned. Depending on the context, we may write $\mathbb{E}_x[\cdot]$ or $\mathbb{E}_{x \sim P}[\cdot]$ to make the random variable and the distribution explicit.

Matrices Throughout, \mathcal{X} will denote a set of interest. Given two subsets $A := (a_1, \dots, a_M)$ and $B := (b_1, \dots, b_N)$ of \mathcal{X} , $k_{AB} \in \mathbb{R}^{M \times N}$ denotes the matrix with elements $[k_{AB}]_{ij} = k(a_i, b_j)$. For a real-valued function $f : \mathcal{X} \rightarrow \mathbb{R}$, $f_A \in \mathbb{R}^M$ denotes the vector with elements $[f_A]_i = f(a_i)$.

Function spaces For a topological space \mathcal{X} , let $C(\mathcal{X})$ denote a set of continuous functions. For a measurable space \mathcal{X} , a measure ν on \mathcal{X} and a constant $p > 0$, let $L_p(\nu)$ be the Banach space of (ν -a.e. equivalent classes of) p -integrable functions with respect to ν :

$$L_p(\nu) := \left\{ f : \mathcal{X} \rightarrow \mathbb{R} : \|f\|_{L_p(\nu)}^p := \int |f(x)|^p d\nu(x) < \infty \right\}. \quad (1)$$

Denote by $L_\nu(\mathcal{X}) := L_2(\nu)$ the one when $\mathcal{X} \subset \mathbb{R}^d$ and ν is the Lebesgue measure. For $f \in L_1(\mathbb{R}^d)$, we define its Fourier transform by

$$\mathcal{F}[f](\omega) := \frac{1}{(2\pi)^{d/2}} \int f(x) e^{-\sqrt{-1} x^\top \omega} dx, \quad \omega \in \mathbb{R}^d.$$

2 Gaussian Processes and RKHSs: Preliminaries

We re-state standard definitions for Gaussian processes (GPs) and RKHSs, reviewing basic properties. Section 2.1 defines positive definite kernels, Sections 2.2 and 2.3 introduce GPs and RKHSs, respectively. Detailed characterizations of GPs and RKHSs can be found in Section 4.

2.1 Positive Definite Kernels

Positive definite kernels play a key role in both Gaussian processes and RKHSs.

Definition 2.1 (Positive definite kernels) Let \mathcal{X} be a nonempty set. A symmetric function $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ is called a positive definite kernel, if for any $n \in \mathbb{N}$, $(c_1, \dots, c_n) \subset \mathbb{R}$ and $(x_1, \dots, x_n) \subset \mathcal{X}$,

$$\sum_{i=1}^n \sum_{j=1}^n c_i c_j k(x_i, x_j) \geq 0.$$

Remark 2.1 Definition 2.1 can be equivalently stated thus: A symmetric function k is positive definite if the matrix $k_{XX} \in \mathbb{R}^{n \times n}$ with elements $[k_{XX}]_{ij} = k(x_i, x_j)$ is positive semidefinite for any finite set $X := (x_1, \dots, x_n) \in \mathcal{X}^n$ of any size $n \in \mathbb{N}$.

In the remainder, for simplicity, *kernel* always means *positive definite kernel*. For $X = (x_1, \dots, x_n) \in \mathcal{X}^n$, the matrix k_{XX} is the *kernel matrix* or *Gram matrix*.

Example 2.1 (Gaussian RBF/Square-Exponential Kernels) Let $\mathcal{X} \subset \mathbb{R}^d$. For $\gamma > 0$, a Gaussian RBF kernel or a square exponential kernel $k_\gamma : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ is defined by

$$k_\gamma(x, x') = \exp\left(-\frac{\|x - x'\|^2}{\gamma^2}\right), \quad x, x' \in \mathcal{X}. \quad (2)$$

In the Gaussian processes literature, to avoid confusion about the term ‘‘Gaussian’’, the kernel (2) is often referred to as the *square-exponential* kernel,¹ while in the kernel literature it is called Gaussian, or Gaussian *radial basis function* (RBF) kernel. The parameter γ determines the length-scale of the associated hypothesis space of functions: As γ increases, the kernel (2) and induced functions change less rapidly, and thus get ‘‘smoother’’ (while they are always infinite differentiable). See Section 4 for details.

Another popular kernel is the Matérn class of functions [Matérn, 1960], which is a standard in the spatial statistics literature [Stein, 1999, Sections 2.7, 2.10]: In fact, Stein [1999, Sec. 1.7] said ‘‘Use the Matérn model’’ as a summary of practical suggestions for modeling spatial data.

Example 2.2 (Matérn kernels) Let $\mathcal{X} \subset \mathbb{R}^d$. For constants $\alpha > 0$ and $h > 0$, the Matérn kernel $k_{\alpha,h} : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ is defined by

$$k_{\alpha,h}(x, x') = \frac{1}{2^{\alpha-1}\Gamma(\alpha)} \left(\frac{\sqrt{2\alpha}\|x - x'\|}{h}\right)^\alpha K_\alpha\left(\frac{\sqrt{2\alpha}\|x - x'\|}{h}\right), \quad x, x' \in \mathcal{X}, \quad (3)$$

where Γ is the gamma function, and K_α is the modified Bessel function of the second kind of order α ,

Remark 2.2 If α can be written as $\alpha = m + 1/2$ for a non-negative integer m , then expression (3) reduces to a product of the exponential function and a polynomial of degree m , which can be computed easily [Rasmussen and Williams, 2006, Section 4.2.1 and Eq. 4.16]:

$$k_{\alpha,h}(x, x') = \exp\left(-\frac{\sqrt{2\alpha}\|x - x'\|}{h}\right) \frac{\Gamma(m+1)}{\Gamma(2m+1)} \sum_{i=1}^m \frac{(m+1)!}{i!(m-i)!} \left(\frac{\sqrt{8\alpha}\|x - x'\|}{h}\right)^{m-i}.$$

¹Sometimes it is also called ‘‘squared exponential’’ or ‘‘exponentiated quadratic.’’

For instance, if $\alpha = 1/2$, $\alpha = 3/2$ or $\alpha = 5/2$, the corresponding Matérn kernels are

$$\begin{aligned} k_{1/2,h}(x, x') &= \exp\left(-\frac{\|x - x'\|}{h}\right), \\ k_{3/2,h}(x, x') &= \left(1 + \frac{\sqrt{3}\|x - x'\|}{h}\right) \exp\left(-\frac{\sqrt{3}\|x - x'\|}{h}\right), \\ k_{5/2,h}(x, x') &= \left(1 + \frac{\sqrt{5}\|x - x'\|}{h} + \frac{5\|x - x'\|^2}{3h^2}\right) \exp\left(-\frac{\sqrt{5}\|x - x'\|}{h}\right). \end{aligned} \tag{4}$$

In particular, (4) is known as the Laplace or exponential kernel.

In the expression (3), the parameter h determines the scale, and α the *smoothness* of functions in the associated hypothesis class: as α increases, the induced functions get smoother. Matérn kernels are appropriate when dealing with “reasonably smooth” (but not very smooth) functions, as the functions in the induced hypothesis class have a finite degree of smoothness [Stein, 1999, Section 6.5]; this is in contrast to a square-exponential kernel, which induces functions with infinite smoothness (i.e., infinite differentiable functions).

Remark 2.3 Square-exponential kernels in Example 2.1 can be obtained as limits of Matérn kernels for $\alpha \rightarrow \infty$ [Stein, 1999, p. 50]. That is, for a Matérn kernel $k_{\alpha,h}$ with $h > 0$ being fixed, we have

$$\lim_{\alpha \rightarrow \infty} k_{\alpha,h}(x, x') = \exp\left(-\frac{\|x - x'\|^2}{2h^2}\right), \quad x, x' \in \mathbb{R}^d.$$

Our last example here is polynomial kernels [Steinwart and Christmann, 2008, Lemma 4.7], which induce hypothesis spaces consisting of polynomials. This class of kernels have been popular in machine learning.

Example 2.3 (Polynomial kernels) Let $\mathcal{X} \subset \mathbb{R}^d$. For $c > 0$ and $m \in \mathbb{N}$, the polynomial kernel $k_{m,c} : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ is defined by

$$k_{m,c}(x, x') = (x^\top x' + c)^m, \quad x, x' \in \mathcal{X}.$$

While we have reviewed here only kernels defined on $\mathcal{X} \subset \mathbb{R}^d$, there are also various kernels defined on non-Euclidian spaces, such as sequences, graphs and distributions; see e.g. Schölkopf and Smola [2002], Schölkopf et al. [2004], Hofmann et al. [2008] and Rasmussen and Williams [2006, Section 4.2]. In fact, as Definition 2.1 indicates, positive definite kernels can be defined on any nonempty set.

2.2 Gaussian Processes

For Gaussian processes, positive definite kernels serve as *covariance functions* of random function values, so they are also called *covariance kernels*. The following definition is taken from Dudley [2002, p. 443].

Definition 2.2 (Gaussian processes) Let \mathcal{X} be a nonempty set, $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ be a positive definite kernel and $m : \mathcal{X} \rightarrow \mathbb{R}$ be any real-valued function. Then a random function $f : \mathcal{X} \rightarrow \mathbb{R}$

is said to be a Gaussian process (GP) with mean function m and covariance kernel k , denoted by $\mathcal{GP}(m, k)$, if the following holds: For any finite set $X = (x_1, \dots, x_n) \subset \mathcal{X}$ of any size $n \in \mathbb{N}$, the random vector

$$\mathbf{f}_X = (f(x_1), \dots, f(x_n))^\top \in \mathbb{R}^n$$

follows the multivariate normal distribution $\mathcal{N}(m_X, k_{XX})$ with covariance matrix $k_{XX} = (k(x_i, x_j))_{i,j=1}^n \in \mathbb{R}^{n \times n}$ and mean vector $m_X = (m(x_1), \dots, m(x_n))^\top$.

Remark 2.4 Definition 2.2 implies that if f is a Gaussian process then there exists a mean function $m : \mathcal{X} \rightarrow \mathbb{R}$ and a covariance kernel $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$. On the other hand, it is also true that for *any* positive definite kernel k and mean function m , there exists a corresponding Gaussian process $f \sim \mathcal{GP}(m, k)$ [Dudley, 2002, Theorem 12.1.3]. There is thus a one-to-one correspondence between Gaussian processes $f \sim \mathcal{GP}(m, k)$ and pairs (m, k) of mean function m and positive definite kernel k .

Remark 2.5 Since k is the covariance function of a Gaussian process, by definition it can be written as

$$k(x, y) = \mathbb{E}_{f \sim \mathcal{GP}(m, k)} [(f(x) - m(x))(f(y) - m(y))], \quad x, y \in \mathcal{X}, \quad (5)$$

where the expectation is with respect to the random function $f \sim \mathcal{GP}(m, k)$. This important property will be used extensively throughout this text.

Remark 2.6 In Definition 2.1, the kernel matrix k_{XX} may be singular: For instance when the kernel k is a polynomial kernel or when some of the points x_1, \dots, x_n are identical. Even in this case, the normal distribution $\mathcal{N}(m_X, k_{XX})$ is well-defined (and thus Definition 2.1 makes sense), while it does not have a density function with respect to the Lebesgue measure; see Dudley [2002, Theorem 9.5.7].

As mentioned in Remark 2.4, the use of a specific kernel k and a mean function m implicitly leads to the use of the corresponding $\mathcal{GP}(k, m)$. Therefore it is practically important to understand the properties of $\mathcal{GP}(k, m)$ (such as smoothness) that are induced by the specification of k and m . For example, if k is a square-exponential kernel on an open set $\mathcal{X} \subset \mathbb{R}^d$, then a sample path $f \sim \mathcal{GP}(0, k)$ is infinitely continuously differentiable, i.e., f is very smooth. In Section 4, we will provide other examples as well as various characterizations for Gaussian processes.

For most practitioners, Gaussian process models manifest themselves in practice much as in their definition, through their finite-dimensional restriction to a concrete set of evaluation nodes; for example a plotting grid. In this case, Gaussian process models are actually very concrete models that are easy to handle on a computer. This fact is sometimes lost in theoretical texts, so we stress it in the following example.

Example 2.4 (GP restricted to finite discrete domain) Let $m : \mathcal{X} \rightarrow \mathbb{R}$ and $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ be a mean and covariance function (kernel), respectively—the arguably most prominent choices are $m(x) \equiv 0$ and $k(x, x') = \exp(-(x - x')^2/2)$. Given a finite set of representer points $X = (x_1, \dots, x_n) \subset \mathcal{X}$, the Algorithm .1 produces a valid draw from the function $f \sim \mathcal{GP}(k, m)$, evaluated at the locations $[f(x_1), \dots, f(x_n)]$. For example, this is how the green draws in Figures 1 and 2 were produced.

Algorithm .1 Concrete algorithm producing independent samples from $\mathcal{GP}(k, m)$ on the finite domain $X \in \mathcal{X}^n$.

```

1 procedure GPSAMPLE( $k, m, X$ )
2    $m_X = [m(x_i)]_{i=1, \dots, n} \in \mathbb{R}^n$  // build mean vector
3    $k_{XX} = [k(x_i, x_j)]_{i, j=1, \dots, n} \in \mathbb{R}^{n \times n}$  // build covariance matrix
4    $R = \text{CHOLESKY}(k_{XX})$  // compute Cholesky decomposition  $k_{XX} = R^T R$ 
5    $u = \text{RANDN}(n, 1)$  // draw  $u \sim \mathcal{N}(0, I_n)$ .
6    $f_X = R^T u + m_X$  // affine transformation of  $u$  gives sample from GP
7 end procedure

```

The following more abstract example, taken from Lindgren et al. [2011], explains that Gaussian processes of Matérn kernels are given as solutions of certain stochastic partial differential equations (SPDE). This was first shown by Whittle [1954, Section 9] for the special case of $d = 2$; see Lindgren et al. [2011] and references therein for further details.

Example 2.5 (GPs of Matérn kernels) Let $k_{\alpha, h}$ be a Matérn kernel in Example 2.2 with parameters $\alpha > 0$ and $h > 0$ defined on $\mathcal{X} = \mathbb{R}$. Then the corresponding Gaussian process $f \sim \mathcal{GP}(0, k_{\alpha, h})$ is the only stationary solution to the following SPDE

$$\left(\frac{2\alpha}{h^2} - \Delta \right)^{(\alpha+d/2)/2} f(x) = \mathbf{w}(x), \quad x \in \mathbb{R}^d,$$

where $\Delta := \sum_{i=1}^d \frac{\partial^2}{\partial x_i^2}$ is the Laplacian and η is the Gaussian white noise process with unit variance.

2.3 Reproducing Kernel Hilbert Spaces

Reproducing kernel Hilbert spaces are defined as follows, where positive definite kernels serve as reproducing kernels.

Definition 2.3 (RKHS) Let \mathcal{X} be a nonempty set and k be a positive definite kernel on \mathcal{X} . A Hilbert space \mathcal{H}_k of functions on \mathcal{X} equipped with an inner-product $\langle \cdot, \cdot \rangle_{\mathcal{H}_k}$ is called a reproducing kernel Hilbert space (RKHS) with reproducing kernel k , if the following are satisfied:

1. For all $x \in \mathcal{X}$, we have $k(\cdot, x) \in \mathcal{H}_k$;
2. For all $x \in \mathcal{X}$ and for all $f \in \mathcal{H}_k$,

$$f(x) = \langle f, k(\cdot, x) \rangle_{\mathcal{H}_k} \quad (\text{Reproducing property}).$$

Remark 2.7 In Definition 2.3, $k(\cdot, x)$ with x being fixed is a real-valued function such that $y \mapsto k(y, x)$ for $y \in \mathcal{X}$. $k(\cdot, x)$ is called the canonical feature map of x , since k can be written as an inner-product in the RKHS as

$$k(x, y) = \langle k(\cdot, x), k(\cdot, y) \rangle_{\mathcal{H}_k}, \quad x, y \in \mathcal{X},$$

which follows from the reproducing property. Therefore $k(\cdot, x)$ is a (possibly infinite dimensional) feature representation of x .

Remark 2.8 By the Moore-Aronszajn theorem [Aronszajn, 1950], for every positive definite kernel k , there exists a unique RKHS \mathcal{H}_k for which k is the reproducing kernel. In this sense, RKHSs and positive definite kernels are one-to-one: for each kernel k there exists a uniquely associated RKHS \mathcal{H}_k , and vice versa.

Given a positive definite kernel k , its RKHS \mathcal{H}_k can be constructed as follows. Let \mathcal{H}_0 be the linear span of feature vectors, that is, each function in \mathcal{H}_0 can be expressed as a finite linear combination of feature vectors:

$$\mathcal{H}_0 := \text{span} \{k(\cdot, x) : x \in \mathcal{X}\} = \left\{ f = \sum_{i=1}^n c_i k(\cdot, x_i) : n \in \mathbb{N}, c_1, \dots, c_n \in \mathbb{R}, x_1, \dots, x_n \in \mathcal{X} \right\}.$$

One can make \mathcal{H}_0 a pre-Hilbert space, by defining an inner-product as follows: For any $f := \sum_{i=1}^n a_i k(\cdot, x_i) \in \mathcal{H}_0$ and $g := \sum_{j=1}^m b_j k(\cdot, y_j) \in \mathcal{H}_0$ with $n, m \in \mathbb{N}$, $a_1, \dots, a_n, b_1, \dots, b_m \in \mathbb{R}$ and $x_1, \dots, x_n, y_1, \dots, y_m \in \mathcal{X}$, the inner-product is defined by

$$\langle f, g \rangle_{\mathcal{H}_0} := \sum_{i=1}^n \sum_{j=1}^m a_i b_j k(x_i, y_j).$$

The norm of \mathcal{H}_0 is induced by the inner-product, i.e., $\|f\|_{\mathcal{H}_0}^2 = \langle f, f \rangle_{\mathcal{H}_0}$. The RKHS \mathcal{H}_k associated with k is then defined as the closure of \mathcal{H}_0 with respect to the norm $\|\cdot\|_{\mathcal{H}_0}$, i.e., $\mathcal{H}_k := \overline{\mathcal{H}_0}$. That is,

$$\mathcal{H}_k = \left\{ f = \sum_{i=1}^{\infty} c_i k(\cdot, x_i) : (c_1, c_2, \dots) \subset \mathbb{R}, (x_1, x_2, \dots) \subset \mathcal{X}, \text{ such that} \right. \quad (6)$$

$$\left. \|f\|_{\mathcal{H}_k}^2 := \lim_{n \rightarrow \infty} \left\| \sum_{i=1}^n c_i k(\cdot, x_i) \right\|_{\mathcal{H}_0}^2 = \sum_{i,j=1}^{\infty} c_i c_j k(x_i, x_j) < \infty \right\}.$$

Remark 2.9 From (6), it is easy to see that functions f in the RKHS \mathcal{H}_k inherit the properties of the kernel k . For instance, if the kernel k is s -times differentiable for $s \in \mathbb{N}$, then so are the functions in \mathcal{H}_k [Steinwart and Christmann, 2008, Corollary 4.36].

An important property of the RKHS norm $\|f\|_{\mathcal{H}_k}$ is that it captures not only the magnitude of a function $f \in \mathcal{H}_k$, but also its *smoothness*: f gets smoother as $\|f\|_{\mathcal{H}_k}$ decreases, and vice versa. This is particularly important in understanding why regularization is required for kernel ridge regression, to avoid overfitting. This smoothness property of the RKHS norm can be seen in the following example on the RKHS of a Matérn kernel, which follows from Rasmussen and Williams [2006, Eq. 4.15] and Wendland [2005, Corollary 10.48]. A complete characterization of RKHSs of Matérn kernels involve Fourier transforms the kernels; see Section 2.4 for details.

Example 2.6 (RKHSs of Matérn kernels: Sobolev spaces) Let $k_{\alpha, h}$ be the Matérn kernel on $\mathcal{X} \subset \mathbb{R}^d$ with Lipschitz boundary² in Example 2.2 with parameters $\alpha > 0$ and $h > 0$

²For the definition of Lipschitz boundary, see e.g., Stein [1970, p.189], Triebel [2006, Definition 4.3] and Kanagawa et al. [2017, Definition 3].

such that $s := \alpha + d/2$ is an integer. Then the RKHS $\mathcal{H}_{k_{\alpha,h}}$ of $k_{\alpha,h}$ is norm-equivalent³ to the Sobolev space $W_2^s(\mathcal{X})$ of order s defined by

$$W_2^s(\mathcal{X}) := \left\{ f \in L_2(\mathcal{X}) : \|f\|_{W_2^s(\mathcal{X})}^2 := \sum_{\beta \in \mathbb{N}_0^d: |\beta| \leq s} \|D^\beta f\|_{L_2(\mathcal{X})}^2 < \infty \right\}. \quad (7)$$

That is, we have $\mathcal{H}_{k_{\alpha,h}} = W_2^s(\mathcal{X})$ as a set of functions, and there exist constants $c_1, c_2 > 0$ such that

$$c_1 \|f\|_{W_2^s(\mathcal{X})} \leq \|f\|_{\mathcal{H}_{k_{\alpha,h}}} \leq c_2 \|f\|_{W_2^s(\mathcal{X})}, \quad \forall f \in \mathcal{H}_{k_{\alpha,h}}, \quad (8)$$

Remark 2.10 The inequality (8) shows the equivalence of the RKHS norm $\|f\|_{\mathcal{H}_{k_{\alpha,h}}}$ and the Sobolev norm $\|f\|_{W_2^s(\mathcal{X})}$ defined in (7). Thus the RKHS norm $\|f\|_{\mathcal{H}_{k_{\alpha,h}}}$ captures the smoothness of the function f with parameter α specifying the order of differentiability. That is, $\|f\|_{\mathcal{H}_{k_{\alpha,h}}}$ takes into account weak derivatives up to order $s = \alpha + d/2$ of the function f . For details of Sobolev spaces, see e.g. Adams and Fournier [2003].

Remark 2.11 The RKHS $\mathcal{H}_{k_{\alpha,h}}$ consists of functions that are weak differentiable up to order $s = \alpha + d/2$. Here one should not confuse the weak differentiability with the classic notion of differentiability. In the classical sense, functions in $\mathcal{H}_{k_{\alpha,h}}$ are only guaranteed to be differentiable up to order α , not $s = \alpha + d/2$; this is a consequence of the Sobolev embedding theorem [Adams and Fournier, 2003, Theorem 4.12]. For definition of weak derivatives, see e.g. Adams and Fournier [2003, Section 1.62]. For instance, consider the case $\alpha = 1/2$, where the kernel is given by (4) and is not differentiable at origin. By Definition 2.3, we have $k_{1/2,h}(\cdot, x) \in \mathcal{H}_{k_{1/2,h}}$ for any $x \in \mathbb{R}^d$; this implies that $\mathcal{H}_{k_{1/2,h}}$ contain functions that are not differentiable in the classical sense.

2.4 A Spectral Characterization for RKHSs Associated with Shift-Invariant Kernels

We provide a characterization of RKHSs associated with shift-invariant kernels on $\mathcal{X} = \mathbb{R}^d$. Recall that a kernel k is called shift-invariant, if it can be written as $k(x, y) = \Phi(x - y)$ for all $x, y \in \mathbb{R}^d$ with a positive definite function $\Phi : \mathbb{R}^d \rightarrow \mathbb{R}$. In the following, the key role is played by the Fourier transform $\mathcal{F}[\Phi]$ of this positive definite function.

Theorem 2.4 below provides a characterization of the RKHS of a shift-invariant kernel in terms of the Fourier transform $\mathcal{F}[\Phi]$ of Φ . This result is available from, e.g., Kimeldorf and Wahba [1970, Lemma 3.1] and Wendland [2005, Theorem 10.12].

Theorem 2.4 *Let k be a shift-invariant kernel on $\mathcal{X} = \mathbb{R}^d$ such that $k(x, y) := \Phi(x - y)$ for $\Phi \in C(\mathbb{R}^d) \cap L_1(\mathbb{R}^d)$. Then the RKHS \mathcal{H}_k of k is given by*

$$\mathcal{H}_k = \left\{ f \in L_2(\mathbb{R}^d) \cap C(\mathbb{R}^d) : \|f\|_{\mathcal{H}_k}^2 = \frac{1}{(2\pi)^{d/2}} \int \frac{|\mathcal{F}[f](\omega)|^2}{\mathcal{F}[\Phi](\omega)} d\omega < \infty \right\}, \quad (9)$$

³Normed vector spaces \mathcal{H}_1 and \mathcal{H}_2 are called norm-equivalent, if $\mathcal{H}_1 = \mathcal{H}_2$ as a set, and if there are constants $c_1, c_2 > 0$ such that $c_1 \|f\|_{\mathcal{H}_2} \leq \|f\|_{\mathcal{H}_1} \leq c_2 \|f\|_{\mathcal{H}_2}$ holds for all $f \in \mathcal{H}_1 = \mathcal{H}_2$, where $\|\cdot\|_{\mathcal{H}_1}$ and $\|\cdot\|_{\mathcal{H}_2}$ denote the norms equipped with \mathcal{H}_1 and \mathcal{H}_2 , respectively.

with the inner-product being

$$\langle f, g \rangle_{\mathcal{H}_k} = \frac{1}{(2\pi)^{d/2}} \int \frac{\mathcal{F}[f](\omega) \overline{\mathcal{F}[g](\omega)}}{\mathcal{F}[\Phi](\omega)} d\omega, \quad f, g \in \mathcal{H}_k,$$

where $\overline{\mathcal{F}[g](\omega)}$ denotes the complex conjugate of $\mathcal{F}[g](\omega)$.

Remark 2.12 Theorem 2.4 shows that the Fourier transform $\mathcal{F}[\Phi]$ determines the members of the RKHS. More specifically, the requirement in (9) shows that, if $\mathcal{F}[\Phi](\omega)$ decays quickly as $|\omega| \rightarrow \infty$, the Fourier transform $\mathcal{F}[f](\omega)$ of each $f \in \mathcal{H}_k$ should also decay quickly as $|\omega| \rightarrow \infty$. Since the tail behaviors of $\mathcal{F}[\Phi]$ and $\mathcal{F}[f]$ determines the smoothness of Φ and f respectively, this implies that if Φ is smooth, f should also be smooth; see examples below.

Remark 2.13 The Fourier transform $\mathcal{F}[\Phi](\omega)$ is known as the *power spectral density* in the stochastic process literature; see e.g. Brémaud [2014, Section 3.3]. It can be written in terms of a certain Fourier transform of the Gaussian process $\mathbf{f} \sim \mathcal{GP}(0, k)$ [Brémaud, 2014, p.161]. We do not explain it in detail, since it requires an explanation of a certain stochastic integral [Brémaud, 2014, Theorem 3.4.1], which is out of the scope of this paper.

The following examples illustrate Theorem 2.4, providing spectral characterizations for RKHSs of square-exponential and Matérn kernels.

Example 2.7 (RKHSs of square-exponential kernels) Let $k_\gamma(x, y) := \Phi_\gamma(x - y) := \exp(-\|x - y\|^2/\gamma^2)$ be the square-exponential kernel with bandwidth $\gamma > 0$ in Example 2.1, and let \mathcal{H}_{k_γ} be the associated RKHS. The Fourier transform of Φ_γ is given by

$$\mathcal{F}[\Phi_\gamma](\omega) = C_{d,\gamma} \exp(-\gamma^2 \|\omega\|^2/4), \quad \omega \in \mathbb{R}^d,$$

where $C_{d,\gamma}$ is a constant depending only on d and γ ; see e.g. Wendland [2005, Theorem 5.20]. Therefore the RKHS \mathcal{H}_{k_γ} can be written as

$$\mathcal{H}_{k_\gamma} = \left\{ f \in L_2(\mathbb{R}^d) \cap C(\mathbb{R}^d) : \|f\|_{\mathcal{H}_{k_\gamma}}^2 = \frac{1}{(2\pi)^{d/2} C_{d,\gamma}} \int |\mathcal{F}[f](\omega)|^2 \exp(\gamma^2 \|\omega\|^2/4) d\omega < \infty \right\},$$

which shows that, for any $f \in \mathcal{H}_{k_\gamma}$, the magnitude of its Fourier transform $|\mathcal{F}[f](\omega)|$ decays exponentially fast as $|\omega| \rightarrow \infty$, and the speed of decay gets quicker as γ increases.

Example 2.8 (RKHSs of Matérn kernels) Let $k_{\alpha,h}$ the Matérn kernel on \mathbb{R}^d with parameters $\alpha > 0$ and $h > 0$ in Example 2.2, and let $\mathcal{H}_{k_{\alpha,h}}$ of $k_{\alpha,h}$ be the associated RKHS. Then $k_{\alpha,h}(x, y) = \Phi_{\alpha,h}(x - y)$ with $\Phi_{\alpha,h}(x) := \frac{2^{1-\alpha}}{\Gamma(\alpha)} (\sqrt{2\alpha} \|x\|/h) K_\alpha(\sqrt{2\alpha} \|x\|/h)$, and the Fourier transform of $\Phi_{\alpha,h}$ is given by

$$\mathcal{F}[\Phi_{\alpha,h}](\omega) = C_{\alpha,h,d} \left(\frac{2\alpha}{h^2} + 4\pi^2 \|\omega\|^2 \right)^{-\alpha-d/2}, \quad \omega \in \mathbb{R}^d, \quad (10)$$

where $C_{\alpha,h,d}$ is a constant depending only on α , h and d ; see, e.g., Rasmussen and Williams [2006, Eq. 4.15]. Therefore the RKHS $\mathcal{H}_{k_{\alpha,h}}$ can be written as

$$\mathcal{H}_{k_{\alpha,h}} = \left\{ f \in L_2(\mathbb{R}^d) \cap C(\mathbb{R}^d) : \|f\|_{\mathcal{H}_{k_{\alpha,h}}}^2 = \frac{1}{(2\pi)^{d/2} C_{\alpha,h,d}} \int |\mathcal{F}[f](\omega)|^2 \left(\frac{2\alpha}{h^2} + 4\pi^2 \|\omega\|^2 \right)^{\alpha+d/2} d\omega < \infty \right\},$$

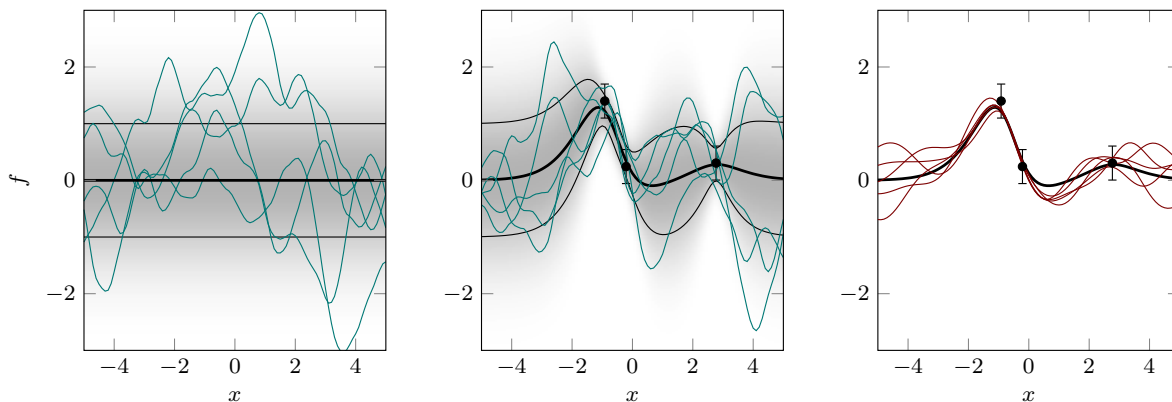


Figure 1: Conceptual sketches of Gaussian process regression (left, center) and kernel ridge regression (right). **Left:** Prior measure $f \sim \mathcal{GP}(0, k)$ with vanishing prior mean and the Matérn-class kernel $k(x, x') = (1 + \sqrt{5}r + 5/3r^2) \exp(-\sqrt{5}r)$ with $r := |x - x'|$. Prior mean function in thick black. Two marginal standard deviations in thin black. Marginal densities as gray shading. 5 samples from prior as green lines. **Center:** Given a dataset (X, Y) of $n = 3$ data points with i.i.d. zero-mean normal noise of standard deviation $\sigma = 0.1$, the posterior measure is also a Gaussian process, with updated mean and covariance functions (all quantities as on the left). **Right:** Kernel ridge regression yields a point estimate (thick black) that is exactly equal to the Gaussian process posterior mean. In contrast to Gaussian process regression, an error estimate is usually not provided. This absence can be deliberate, as one may not be willing to impose the assumptions necessary to define such an estimate (e.g., additive Gaussian noise assumption). For comparison with the GP samples, the plot also shows some functions with the property that $f_X^\top k_{XX}^{-1} f_X = \|f_X\|_{k_{XX}^{-1}}^2 = 1$ (but KRR does not assume the true function is of this form).

which shows that, for any $f \in \mathcal{H}_{k_{\alpha, h}}$, the magnitude of its Fourier transform $|\mathcal{F}[f](\omega)|$ decays polynomially fast as $|\omega| \rightarrow \infty$, and the speed of decay gets quicker as α increases. Moreover, from (10) and Wendland [2005, Corollary 10.48], it follows that $\mathcal{H}_{k_{\alpha, h}}$ is norm-equivalent to the Sobolev space of order $\alpha + d/2$.

3 Connections between Gaussian Process and Kernel Ridge Regression

Regression is a fundamental task in statistics and machine learning. The *interpolation* problem is regression with noise-free observations and has been studied mainly in the literature on numerical analysis, and more recently in the context of Bayesian optimization. We compare two approaches to these problems based on Gaussian processes and kernel methods, namely *Gaussian process regression* and *kernel ridge regression* (see also Figure 1 for illustration).

We first describe the problem of regression, and set notation. Let \mathcal{X} be a nonempty set and $f : \mathcal{X} \rightarrow \mathbb{R}$ be a function. Assume that one is given a set of pairs $(x_i, y_i)_{i=1}^n \subset \mathcal{X} \times \mathbb{R}$ for

$n \in \mathbb{N}$, which is referred to as *training data*, such that

$$y_i = \mathbf{f}(x_i) + \xi_i, \quad i = 1, \dots, n, \quad (11)$$

where ξ_i is a zero-mean random variable that represents “noise” in the output, or the variability in the responses which is not explained by the input vectors. The task of regression is to estimate the unknown function \mathbf{f} based on the training data $(x_i, y_i)_{i=1}^n$. The function \mathbf{f} is called the *regression function*, and is the conditional expectation of the output given an input:

$$\mathbf{f}(x) = \mathbb{E}[y|x],$$

where (x, y) is a random variable with the conditional distribution of y given x following the model (11).

If there is no output noise, i.e., $\xi_i = 0$, the problem is called *interpolation*; in this case one can obtain exact function values $y_i = \mathbf{f}(x_i)$ for training. We will frequently use the notation $X := (x_1, \dots, x_n) \in \mathcal{X}^n$ for the set of input data points, and $Y := (y_1, \dots, y_n)^\top \in \mathbb{R}^n$ for the set of outputs (or $\mathbf{f}_X := (\mathbf{f}(x_1), \dots, \mathbf{f}(x_n))^\top \in \mathbb{R}^n$ in the noise free case).

This section first reviews Gaussian process regression and interpolation in Section 3.1, and kernel ridge regression and kernel interpolation in Section 3.2. We summarize and discuss equivalences between the two approaches in Section 3.3. In GP-regression, the *posterior variance* function plays a fundamental role, but its kernel interpretation has not been well understood. In Section 3.4, we show that there exists an interpretation of the posterior variance function as a certain *worst case error* in the RKHS. Coming back to regression itself, in Section 3.5 we provide a weight-vector viewpoint for the regression problem, and discuss the equivalence between regularization and an additive noise assumption.

3.1 Gaussian Process Regression and Interpolation

Gaussian process regression, also known as *Kriging* or *Wiener-Kolmogorov prediction*, is a Bayesian nonparametric method for regression. Being a Bayesian approach, GP-regression produces a *posterior distribution* of the unknown regression function \mathbf{f} , provided the training data (X, Y) , a prior distribution Π_0 on \mathbf{f} , and a likelihood function denoted by $\ell_{X,Y}(\mathbf{f})$. More specifically, the prior Π_0 is defined as a Gaussian process $\mathcal{GP}(m, k)$ with mean function $m : \mathcal{X} \rightarrow \mathbb{R}$ and covariance kernel $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$, i.e.,

$$\mathbf{f} \sim \mathcal{GP}(m, k). \quad (12)$$

Since this GP serves as a prior, the mean function m and the kernel k should be chosen so that they reflect one’s prior knowledge or belief about the regression function \mathbf{f} ; this will be discussed later.

On the other hand, a likelihood function is defined by a probabilistic model $p(y_i|\mathbf{f}(x_i))$ for the noise variables ξ_1, \dots, ξ_n , since this determines the distribution of the observations $Y = (y_1, \dots, y_n)^\top$ with the additive noise model (11). It is typical to assume that ξ_1, \dots, ξ_n are i.i.d. centered Gaussian random variables with variance $\sigma^2 > 0$:

$$\xi_i \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \sigma^2), \quad i = 1, \dots, n. \quad (13)$$

Thus the likelihood function is defined as

$$\ell_{X,Y}(\mathbf{f}) = \prod_{i=1}^n \mathcal{N}(y_i|\mathbf{f}(x_i), \sigma^2), \quad (14)$$

where $\mathcal{N}(\cdot|\mu, \sigma^2)$ denotes the density function of the normal distribution $\mathcal{N}(\mu, \sigma^2)$ of mean μ and variance σ^2 . In general however, GP-regression allows the noise variables to be correlated Gaussian with varying magnitudes of variances.

By Bayes' rule, the posterior distribution $\Pi_n(\mathfrak{f}|Y, X)$ is then given as

$$d\Pi_n(\mathfrak{f}|X, Y) \propto \ell_{X,Y}(\mathfrak{f})d\Pi_0(\mathfrak{f}) = \prod_{i=1}^n \mathcal{N}(y_i|\mathfrak{f}(x_i), \sigma^2)d\Pi_0(\mathfrak{f}). \quad (15)$$

As shown in the following theorem, which is well known in the literature, the posterior $\Pi_n(\mathfrak{f}|X, Y)$ is again a Gaussian process, whose mean function and covariance function are obtained by simple linear algebra.

Theorem 3.1 *Assume (11), (12) and (13), and let $X = (x_1, \dots, x_n) \in \mathcal{X}^n$ and $Y = (y_1, \dots, y_n)^\top \in \mathbb{R}^n$. Then we have*

$$\mathfrak{f}|Y \sim \mathcal{GP}(\bar{m}, \bar{k}),$$

where $\bar{m} : \mathcal{X} \rightarrow \mathbb{R}$ and $\bar{k} : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ are given by

$$\bar{m}(x) = m(x) + k_{xX}(k_{XX} + \sigma^2 I_n)^{-1}(Y - m_X), \quad x \in \mathcal{X}, \quad (16)$$

$$\bar{k}(x, x') = k(x, x') - k_{xX}(k_{XX} + \sigma^2 I_n)^{-1}k_{Xx'}, \quad x, x' \in \mathcal{X}, \quad (17)$$

where $k_{Xx} = k_{xX}^\top = (k(x_1, x), \dots, k(x_n, x))^\top$.

As $\mathcal{GP}(\bar{m}, \bar{k})$ is a posterior Gaussian process, \bar{m} is referred to as the *posterior mean function* and \bar{k} the *posterior covariance function*. It is instructive to see how the Gaussian noise assumptions (13) and the GP prior (12) lead to the closed form expressions (16) and (17), because this can be done *without relying on Bayes' rule*. This is important for the following two reasons: (i) Since the prior and posterior are defined on an *infinite* dimensional space of functions, Bayes' rule is more involved and thus does not produce the expressions (16) and (17) directly (see e.g. Stuart 2010, Theorem 6.31); (ii) When dealing with the *noise-free* setting where $\sigma^2 = 0$, Bayes' rule cannot be used because the likelihood function is degenerate [Cockayne et al., 2017].

To prove Theorem 3.1, first recall a basic formula for conditional distributions of Gaussian random vectors (see e.g. Rasmussen and Williams 2006, Appendix A.2).

Proposition 3.2 *Let $a \in \mathbb{R}^n$ and $b \in \mathbb{R}^m$ be Gaussian random vectors such that*

$$\begin{bmatrix} a \\ b \end{bmatrix} \sim \mathcal{N} \left(\begin{bmatrix} \mu_a \\ \mu_b \end{bmatrix}, \begin{bmatrix} A & C \\ C^\top & B \end{bmatrix} \right), \quad (18)$$

where $\mu_a \in \mathbb{R}^n$, $\mu_b \in \mathbb{R}^m$ are the mean vectors, $A \in \mathbb{R}^{n \times n}$, $B \in \mathbb{R}^{m \times m}$ are the covariance matrices (where A is strictly positive definite), and $C \in \mathbb{R}^{n \times m}$. Then the conditional distribution of b given a is

$$b|a \sim \mathcal{N}(\mu_b + C^\top A^{-1}(a - \mu_a), B - C^\top A^{-1}C). \quad (19)$$

Proof [Theorem 3.1] Let $m \in \mathbb{N}$, and let $Z = (z_1, \dots, z_m) \in \mathcal{X}^m$ be any finite set of points. Then the observations $Y \in \mathbb{R}^n$ and GP-function values $\mathfrak{f}_Z = (\mathfrak{f}(z_1), \dots, \mathfrak{f}(z_m))^\top \in \mathbb{R}^m$ are jointly Gaussian random vectors such that

$$\begin{bmatrix} Y \\ \mathfrak{f}_Z \end{bmatrix} \sim \mathcal{N} \left(\begin{bmatrix} m_X \\ m_Z \end{bmatrix}, \begin{pmatrix} k_{XX} + \sigma^2 I_n & k_{XZ} \\ k_{ZX} & k_{ZZ} \end{pmatrix} \right).$$

In the notation of (18), this corresponds to $a = Y|X$, $b = \mathbf{f}_Z$, $\mu_a = m_X$, $\mu_b = m_Z$, $A = k_{XX} + \sigma^2 I_n$, $B = k_{ZZ}$ and $C = k_{XZ}$. Applying the formula (19) in Proposition 3.2, the conditional distribution of \mathbf{f}_Z given Y is then given as

$$\mathbf{f}_Z|Y \sim \mathcal{N}(\bar{\mu}, \bar{\Sigma}),$$

where

$$\begin{aligned} \bar{\mu} &:= m_Z + k_{ZX}(k_{XX} + \sigma^2 I_n)^{-1}(Y - m_X) \in \mathbb{R}^m, \\ \bar{\Sigma} &:= k_{ZZ} - k_{ZX}(k_{XX} + \sigma^2 I_n)^{-1}k_{XZ} \in \mathbb{R}^{m \times m}. \end{aligned}$$

This mean vector and the covariance matrix can be written as $\bar{\mu} := \bar{m}_Z$, $\bar{\Sigma} = \bar{k}_{ZZ}$, where $\bar{m} : \mathcal{X} \rightarrow \mathbb{R}$ and $\bar{k} : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ are defined as (16) and (17). In other words,

$$\mathbf{f}_Z|Y \sim \mathcal{N}(\bar{m}_Z, \bar{k}_{ZZ}). \quad (20)$$

Note that (20) holds for any set of points $Z = (z_1, \dots, z_m) \in \mathcal{X}^m$ of any size $m \in \mathbb{N}$. Therefore, by the Kolmogorov extension theorem [Dudley, 2002, Theorems 12.1.2] and the definition of GPs (Definition 2.2), this implies that the process $\mathbf{f} \sim \mathcal{GP}(m, k)$ conditioned on the training data X, Y is a draw from $\mathcal{GP}(\bar{m}, \bar{k})$. ■

Remark 3.1 Given a test input x , prediction of the output value $\mathbf{f}(x)$ is carried out by evaluating the posterior mean function (16), as we have $\mathbb{E}[\mathbf{f}(x)|X, Y] = \bar{m}(x)$ by definition. On the other hand, the posterior covariance \bar{k} can be used to quantify uncertainties over output values; this will be discussed in Section 3.4.

Remark 3.2 As it can be seen from the expressions (16) and (17), \bar{m} and \bar{k} depend on the choice of the prior mean function m , the kernel k and the noise variance σ^2 . These are hyper-parameters of GP-regression, and the determination of them can be carried out, for example, by the empirical Bayes method, i.e., maximization of the marginal likelihood of the data given hyperparameters (for regression this is available in closed form); see Rasmussen and Williams [2006] for details.

Noise-free case: Gaussian process interpolation. Consider the noise-free case where exact function values $y_i = \mathbf{f}(x_i)$, $i = 1, \dots, n$ are provided for training. In this case, the likelihood function (14) is degenerate and thus not well-defined, since the distribution of y_i given $\mathbf{f}(x_i)$ is the Dirac distribution at $\mathbf{f}(x_i)$, which has no density function. Thus, it is not possible to apply Bayes' rule to derive the posterior distribution of \mathbf{f} as in (15); see also Cockayne et al. [2017, Section 2.5]. However, as the proof for Theorem 3.1 indicates, the conditional distribution of \mathbf{f} given training data $(x_i, \mathbf{f}(x_i))_{i=1}^n$ can be derived based on Gaussian calculus, without relying on Bayes' rule. The resulting posterior mean function and covariance function are respectively given as (16) and (17) with $\sigma^2 = 0$, as shown in the following theorem.

Theorem 3.3 *Assume (12), and let $X = (x_1, \dots, x_n) \in \mathcal{X}^n$ and $\mathbf{f}_X = (\mathbf{f}(x_1), \dots, \mathbf{f}(x_n))^\top \in \mathbb{R}^n$. Moreover, assume that the kernel matrix $k_{XX} = (k(x_i, x_j))_{i,j=1}^n \in \mathbb{R}^{n \times n}$ is invertible. Then the conditional distribution of \mathbf{f} given (X, \mathbf{f}_X) is a Gaussian process*

$$\mathbf{f} | \mathbf{f}_X \sim \mathcal{GP}(\bar{m}, \bar{k}),$$

where $\bar{m} : \mathcal{X} \rightarrow \mathbb{R}$ and $\bar{k} : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ are given by

$$\bar{m}(x) = m(x) + k_{xX} k_{XX}^{-1} (\mathbf{f}_X - m_X), \quad x \in \mathcal{X}, \quad (21)$$

$$\bar{k}(x, x') = k(x, x') - k_{xX} k_{XX}^{-1} k_{Xx'}, \quad x, x' \in \mathcal{X}. \quad (22)$$

Proof Since k_{XX} is assumed to be invertible, the assertion can be proven by modifying the proof of Theorem 3.1. Specifically, this can be done by replacing $k_{XX} + \sigma^2 I_n$ in the proof of Theorem 3.1 by k_{XX} , and Y by \mathbf{f}_X . ■

Remark 3.3 In Theorem 3.3, the kernel matrix k_{XX} is required to be invertible. If this condition is not satisfied, then the expressions (21) are (22) not well-defined. For instance, k_{XX} is not invertible, if some of the points in $X = (x_1, \dots, x_n)$ are identical, or if the kernel k is a polynomial kernel of order m such that $n > m$.

This way of using Gaussian processes in modeling *deterministic* functions is becoming popular in machine learning, in particular in the context of Bayesian optimization (e.g., Bull, 2011) as well as in the emerging field of probabilistic numerics [Hennig et al., 2015]: For instance, Bayesian quadrature, a probabilistic numerics approach to numerical integration, involves integration of a fixed deterministic function, which is modeled as a Gaussian process with noise-free outputs; see Section 6.2 for details.

The noise-free situation appears for instance when a measurement equipment for the output values is very accurate, or when the function values are obtained as a result of computer experiments. In the latter case, GP-interpolation is often called *emulation* in the literature. In these situations, typically the function of interest is very expensive to evaluate, so inference should be done based on a small number of function evaluations. Gaussian processes are useful for this purpose, since one can gain statistical efficiency by incorporating available prior knowledge about the function via the choice of a covariance kernel.

Remark 3.4 For the noise-free case, a posterior distribution may be well-defined by assuming the existence of very small noise in outputs, which corresponds to applying regularization with a very small regularization constant; this is called “jitter” in the kriging literature. This is practically reasonable, since if the kernel matrix k_{XX} is singular (or close to singular, leading to numerical issues), then the posterior mean (21) as well as the posterior variance (22) are not well-defined without regularization.

3.2 Kernel Ridge Regression and Kernel Interpolation

Kernel ridge regression (KRR), which is also known as *regularized least-squares* [Caponnetto and Vito, 2007] or *spline smoothing* [Wahba, 1990], arises as a regularized empirical risk minimization problem where the hypothesis space is chosen to be an RKHS \mathcal{H}_k . That is, we are interested in solving the problem

$$\hat{f} = \arg \min_{f \in \mathcal{H}_k} \frac{1}{n} \sum_{i=1}^n L(x_i, y_i, f(x_i)) + \lambda \|f\|_{\mathcal{H}_k}^2.$$

where $L : \mathcal{X} \times \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}^+$ is a loss function, and $\lambda > 0$ is a regularization constant. The loss function penalizes the deviations between predicted outputs $f(x_i)$ and true outputs y_i . The

regularization constant λ controls the smoothness of the estimator, to avoid overfitting: the larger the λ is, the smoother the resulting estimator \hat{f} becomes. Regularization is necessary, as nonparametric estimation of a function from a finite sample is an ill-posed inverse problem, given also that output values are contaminated by noise.

The KRR estimator then arises when using the square loss $L(x, y, y') = (y - y')^2$:

$$\hat{f} = \arg \min_{f \in \mathcal{H}_k} \frac{1}{n} \sum_{i=1}^n (f(x_i) - y_i)^2 + \lambda \|f\|_{\mathcal{H}_k}^2. \quad (23)$$

While this least-square problem is over the function space \mathcal{H}_k , which may be infinite dimensional, its solution can be obtained by simple linear algebra, as the following theorem shows. As it is simple and instructive, we show its proof based on the representer theorem [Schölkopf et al., 2001].

Theorem 3.4 *If $\lambda > 0$, the solution to (23) is unique as a function, and is given by*

$$\hat{f}(x) = k_{xX}(k_{XX} + n\lambda I_n)^{-1}Y = \sum_{i=1}^n \alpha_i k(x, x_i), \quad x \in \mathcal{X}, \quad (24)$$

where

$$(\alpha_1, \dots, \alpha_n)^\top := (k_{XX} + n\lambda I_n)^{-1}Y \in \mathbb{R}^n. \quad (25)$$

If we further assume that k_{XX} is invertible, then the coefficients $(\alpha_1, \dots, \alpha_n)$ in (24) are uniquely given by (25).

Proof Because of the regularization term in (23), one can apply the representer theorem [Schölkopf et al., 2001, Theorem 1]. This implies that the solution to (23) can be written as a weighted sum of feature vectors $k(\cdot, x_1), \dots, k(\cdot, x_n)$, i.e.,

$$\hat{f} = \sum_{i=1}^n \alpha_i k(\cdot, x_i), \quad (26)$$

for some coefficients $\alpha_1, \dots, \alpha_n \in \mathbb{R}$. Let $\boldsymbol{\alpha} := (\alpha_1, \dots, \alpha_n)^\top \in \mathbb{R}^n$. By substituting the expression (26) in (23), the optimization problem now becomes

$$\min_{\boldsymbol{\alpha} \in \mathbb{R}^n} \frac{1}{n} [\boldsymbol{\alpha}^\top k_{XX}^2 \boldsymbol{\alpha} - 2\boldsymbol{\alpha}^\top k_{XX} Y + \|Y\|^2] + \lambda \boldsymbol{\alpha}^\top k_{XX} \boldsymbol{\alpha}, \quad (27)$$

where we used $\boldsymbol{\alpha}^\top k_{XX} \boldsymbol{\alpha} = \|\hat{f}\|_{\mathcal{H}_k}^2$, which follows from the reproducing property. Differentiating this objective function with respect to $\boldsymbol{\alpha}$, setting it equal to 0 and arranging the resulting equation yields

$$k_{XX}(k_{XX} + n\lambda I_n)\boldsymbol{\alpha} = k_{XX}Y. \quad (28)$$

Obviously $\boldsymbol{\alpha} = (k_{XX} + n\lambda I_n)^{-1}Y$ is one of the solutions to (28). Since the objective function in (27) is a convex function of $\boldsymbol{\alpha}$ (while it may not be strictly convex unless k_{XX} is strictly positive definite or invertible), $\boldsymbol{\alpha}$ attains the minimum of the objective function. Since the objective function in (27) is equal to that of (23), the function (26) with $\boldsymbol{\alpha} = (k_{XX} + n\lambda I_n)^{-1}Y$ attains the minimum of (23).

Note that since the square loss is convex⁴, the solution to (23) is unique as a function [Steinwart and Christmann, 2008, Theorem 5.5]. Hence (26) with $\boldsymbol{\alpha} = (k_{XX} + n\lambda I_n)^{-1}Y$ gives the unique solution to (23) as a function, and this proves the first claim. If k_{XX} is further invertible, (28) reduces to $(k_{XX} + n\lambda I_n)\boldsymbol{\alpha} = Y$, from which the second claim follows. \blacksquare

Remark 3.5 While Theorem 3.4 shows that (24) is the unique solution of (23) as a *function*, this does not mean that the *coefficients* $\alpha_1, \dots, \alpha_n$ in (24) are uniquely determined, unless the kernel matrix k_{XX} is invertible. This is because there may be multiple solutions to the linear system (28), if k_{XX} is not invertible. (More precisely, if $(k_{XX} + n\lambda I_n)\boldsymbol{\alpha}' - Y$ is in the null space of k_{XX} , such an $\boldsymbol{\alpha}'$ is a solution to (28), even when $(k_{XX} + n\lambda I_n)\boldsymbol{\alpha}' - Y \neq 0$.) However, even when multiple solutions to (28) exist, they result in the same estimator (24) as a function, which can be shown as follows. Therefore one can always use the coefficients given in (25).

Let $\boldsymbol{\alpha}' := (\alpha'_1, \dots, \alpha'_n)^\top \in \mathbb{R}^n$ be another solution to (28). As mentioned in the proof, since the objective function (27) is a convex function of $\boldsymbol{\alpha}$, this solution $\boldsymbol{\alpha}'$ also attains the minimum of the objective function in (27), and thus the resulting function $\hat{f}' := \sum_{i=1}^n \alpha'_i k(\cdot, x_i)$ attains the minimum of the objective function in (23). However, since the solution to (23) is unique [Steinwart and Christmann, 2008, Theorem 5.5], we have $\hat{f} = \hat{f}'$, where \hat{f} is the KRR estimator (24).

Noise-free case: Kernel interpolation In the noise-free case where $y_i = f(x_i)$, $i = 1, \dots, n$, the estimator of f is given by (24) with $\lambda = 0$; that is,

$$\hat{f}(x) = k_{xX} k_{XX}^{-1} f_X = \sum_{i=1}^n \alpha_i k(x, x_i), \quad x \in \mathcal{X}, \quad (29)$$

where $f_X := (f(x_1), \dots, f(x_n))^\top \in \mathbb{R}^n$ and $(\alpha_1, \dots, \alpha_n)^\top := k_{XX}^{-1} f_X$. Thus, in this case the kernel matrix k_{XX} is required to be invertible. The estimator (29) is obtained as a solution of the following optimization problem in the RKHS. We provide a proof based on that of Berlinet and Thomas-Agnan [2004, Theorem 58 in p. 112].

Theorem 3.5 *Let k be a kernel on a nonempty set \mathcal{X} , and $X = (x_1, \dots, x_n) \in \mathcal{X}^n$ be such that the kernel matrix k_{XX} is invertible. Then (29) is the unique solution of the following optimization problem:*

$$\hat{f} := \arg \min_{f \in \mathcal{H}_k} \|f\|_{\mathcal{H}_k} \quad \text{subject to} \quad f(x_i) = f(x_i), \quad i = 1, \dots, n. \quad (30)$$

Proof Let \mathcal{S}_0 be the linear span of the feature vectors $k(\cdot, x_1), \dots, k(\cdot, x_n)$, that is,

$$\mathcal{S}_0 := \left\{ f = \sum_{i=1}^n \alpha_i k(\cdot, x_i) : \alpha_1, \dots, \alpha_n \in \mathbb{R} \right\}.$$

⁴A loss function $L : \mathcal{X} \times \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$ is called convex, if $L(x, y, \cdot) : \mathbb{R} \rightarrow \mathbb{R}$ is convex for all fixed $x \in \mathcal{X}$ and $y \in \mathbb{R}$ [Steinwart and Christmann, 2008, Definition 2.12]

Let \mathcal{H}_0 be the set of all functions in \mathcal{H}_k that interpolate the data $(x_i, \mathfrak{f}(x_i))_{i=1}^n$:

$$\mathcal{H}_0 := \{f \in \mathcal{H}_k : f(x_i) = \mathfrak{f}(x_i), \quad i = 1, \dots, n\}.$$

It is easy to see that (29) satisfies $\hat{f}(x_\ell) = \mathfrak{f}(x_\ell)$ for all $\ell = 1, \dots, n$, and thus $\hat{f} \in \mathcal{S}_0 \cap \mathcal{H}_0$.

We first show that $\mathcal{S}_0 \cap \mathcal{H}_0$ consists only of \hat{f} . To this end, assume that there exist another $g \in \mathcal{S}_0 \cap \mathcal{H}_0$, and let $g = \sum_{i=1}^n \beta_i k(\cdot, x_i)$ for $\beta_1, \dots, \beta_n \in \mathbb{R}$. Then

$$\hat{f} - g = \sum_{i=1}^n (\alpha_i - \beta_i) k(\cdot, x_i) \in \mathcal{S}_0.$$

On the other hand, since $\hat{f}(x_\ell) = g(x_\ell) = \mathfrak{f}(x_\ell)$ for all $\ell = 1, \dots, n$, we have

$$\hat{f}(x_\ell) - g(x_\ell) = \left\langle \hat{f} - g, k(\cdot, x_\ell) \right\rangle_{\mathcal{H}_k} = 0, \quad \forall \ell = 1, \dots, n.$$

This implies that $\hat{f} - g \in \mathcal{S}_0^\perp$, where $\mathcal{S}_0^\perp \subset \mathcal{H}_k$ is the orthogonal complement of \mathcal{S}_0 . Therefore $\hat{f} - g \in \mathcal{S}_0 \cap \mathcal{S}_0^\perp = \{0\}$, which implies that $\hat{f} = g$. Thus, $\mathcal{S}_0 \cap \mathcal{H}_0 = \{\hat{f}\}$.

Finally, we show that \bar{f} is the solution of (30). It is easy to show that \mathcal{H}_0 is convex and closed. Thus there exists an element $f^* \in \mathcal{H}_0$ such that

$$f^* = \arg \min_{f \in \mathcal{H}_0} \|f\|_{\mathcal{H}_k}.$$

For any $v \in \mathcal{S}_0^\perp$, we have $\langle f^* + v, k(\cdot, x_\ell) \rangle_{\mathcal{H}_k} = \langle f^*, k(\cdot, x_\ell) \rangle_{\mathcal{H}_k} = f^*(x_\ell) = \mathfrak{f}(x_\ell)$ for all $\ell = 1, \dots, n$, and thus $f^* + v \in \mathcal{H}_0$. By definition, $\|f^*\|_{\mathcal{H}_k} \leq \|f^* + v\|_{\mathcal{H}_k}$, and this holds for all $v \in \mathcal{S}_0^\perp$. This implies that f^* belongs to the orthogonal complement of \mathcal{S}_0^\perp , which is \mathcal{S}_0 since \mathcal{S}_0 is closed. That said, $f^* \in \mathcal{S}_0$ and thus $f^* \in \mathcal{H}_0 \cap \mathcal{S}_0 = \{\hat{f}\}$, which implies $f^* = \hat{f}$. ■

Remark 3.6 As the RKHS norm $\|f\|_{\mathcal{H}_k}$ quantifies the smoothness of the function $f \in \mathcal{H}_k$ (see Example 2.6 and Section 4 for this property of the RKHS norm), the solution to the optimization problem (30) is interpreted as the smoothest function in the RKHS that passes all the training data $(x_1, \mathfrak{f}(x_1)), \dots, (x_n, \mathfrak{f}(x_n))$.

Remark 3.7 In practice, regularization for matrix inversion in (29) may still be needed even if there is no output noise, for the sake of numerical stability when the kernel matrix is nearly singular. For this purpose, nevertheless, the regularization constant may be chosen to be very small, since it is not relevant to the variance of output noise. See Wendland and Rieger [2005, Section 3.4] and Schaback and Wendland [2006, Section 7.8] for theoretical supports.

3.3 Equivalences in Regression and Interpolation

From the expressions (16) and (24), it is immediate that the following equivalence holds for GP-regression and kernel ridge regression. While this result has been well known in the literature, we summarize it in the following proposition.

Proposition 3.6 *Let k be a positive definite kernel on a nonempty set \mathcal{X} and $(x_i, y_i)_{i=1}^n \subset \mathcal{X} \times \mathbb{R}$ be training data, and define $X := (x_1, \dots, x_n) \in \mathcal{X}^n$ and $Y := (y_1, \dots, y_n)^\top \in \mathbb{R}^n$. Then we have $\bar{m} = \hat{f}$ if $\sigma^2 = n\lambda$, where*

- \bar{m} is the posterior mean function (16) of GP-regression based on (X, Y) , the GP prior $\mathfrak{f} \sim \mathcal{GP}(0, k)$ and the modeling assumption (11), where $\xi_1, \dots, \xi_n \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \sigma^2)$ with variance $\sigma^2 > 0$;
- \hat{f} is the solution (24) to kernel ridge regression (23) based on (X, Y) , the RKHS \mathcal{H}_k , and regularization constant $\lambda > 0$.

Remark 3.8 One immediate consequence of Proposition 3.6 is that the posterior mean function \bar{m} belongs to the RKHS \mathcal{H}_k , under the assumptions in Proposition 3.6. On the other hand, it is well known that a sample $\mathfrak{f} \sim \mathcal{GP}(\bar{m}, \bar{k})$ from the posterior GP does not belong to \mathcal{H}_k almost surely; we will discuss why this is the case in Section 4, and see nevertheless that the GP sample belongs to a certain RKHS induced by \mathcal{H}_k , which is larger than \mathcal{H}_k .

Remark 3.9 Proposition 3.6 implies that the additive Gaussian noise assumption (11) in GP regression plays the role of regularization in KRR, as the two estimators are identical if $\lambda = \sigma^2/n$. Recall that λ controls the smoothness of \hat{f} : as λ increases, \hat{f} gets smoother. Therefore, the assumption that the noise variance σ^2 is large amounts to the assumption that the latent function \mathfrak{f} is smoother than the observed process. This interpretation may be explained in the following way. Let η be the zero-mean Gaussian process with a covariance kernel $\delta : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ defined by

$$\delta(x, x') = \begin{cases} 1 & (x = x') \\ 0 & (x \neq x'). \end{cases} \quad (31)$$

Note that this is a valid kernel, since it is positive definite. Then the noise variable ξ_i can be written as $\xi_i = \sigma\eta(x_i)$, since this results in $\xi_1, \dots, \xi_n \stackrel{i.i.d.}{\sim} \mathcal{N}(0, 1)$. Define a Gaussian process y by

$$y = \mathfrak{f} + \sigma\eta. \quad (32)$$

where $\mathfrak{f} : \mathcal{X} \rightarrow \mathbb{R}$ is the latent function. The training observations y_i can then be given as evaluations of the process (32), that is $y_i = y(x_i)$, $i = 1, \dots, n$. Thus, the problem of regression is to infer the latent function \mathfrak{f} based on evaluations of the process (32), i.e., $(x_i, y(x_i))_{i=1}^n$. Knowing the model (32), which states that y is a noisy version of \mathfrak{f} , one knows that the observed process y must be rougher than the latent function \mathfrak{f} , or that \mathfrak{f} must be smoother than y . In other words, assuming the noise model (32) amounts to assuming the latent function \mathfrak{f} being smoother than the observed process y ; this is how the noise assumption plays the role of regularization.

Noise-free case: interpolation. For the noise-free case, there also exists an equivalence between GP-interpolation and kernel interpolation, which we summarize in the following result.

Proposition 3.7 *Let k be a positive definite kernel on a nonempty set \mathcal{X} and $(x_i, \mathfrak{f}(x_i))_{i=1}^n \subset \mathcal{X} \times \mathbb{R}$ be training data, and define $X := (x_1, \dots, x_n) \in \mathcal{X}^n$ and $\mathfrak{f}_X := (\mathfrak{f}(x_1), \dots, \mathfrak{f}(x_n))^T \in \mathbb{R}^n$. Assume that the kernel matrix k_{XX} is invertible. Then we have $\bar{m} = \hat{f}$, where*

- \bar{m} is the posterior mean function (21) of GP-interpolation based on (X, \mathfrak{f}_X) and the GP prior $\mathfrak{f} \sim \mathcal{GP}(0, k)$;
- \hat{f} is the solution (29) to kernel interpolation (30) based on (X, \mathfrak{f}_X) and the RKHS \mathcal{H}_k .

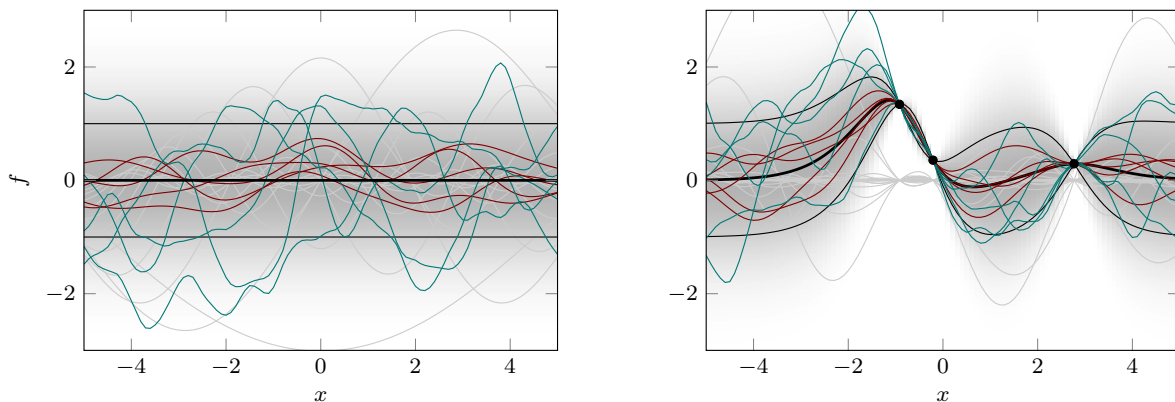


Figure 2: In-model error estimation. Plots similar to Fig. 1. **Left:** Hypothesis class/prior: The plot shows five sample paths from the GP prior in green and, for comparison, five functions with $f_X^\top k_{XX}^{-1} f_X = 1$ in red. In light gray in the background: Eigenfunction spectrum (regular grid over the continuous space of such functions), scaled by their eigenvalues. (See Sec. 4.1.1 for eigen expansions of GP and RKHSs) **Right:** When constrained on noise-less observations, both Gaussian process regression and kernel ridge regression afford the same in-model error estimate, plotted as two thin black lines (Proposition 3.10). In the GP context, this is the error bar of one marginal standard deviation. In the kernel context, it is the worst case error if the true function has unit RKHS norm. The red functions (which approximate such unit-norm RKHS elements) lie entirely inside this region, while GP samples (green) lie inside it for $\sim 68\%$ of the path (the expected value, the Gaussian probability mass within one standard-deviation). Another visible feature is that the GP samples are rougher than the unit-norm representers.

3.4 Error Estimates: Posterior Variance and Worst-Case Error

Gaussian process regression is usually employed in settings that also call for a notion of *uncertainty*, or *error estimate*. The object given this interpretation is the posterior covariance function \bar{k} given by (17) or its scalar value $\bar{k}(x, x)$ at a particular location $x \in \mathcal{X}$; this is the (*marginal*) *posterior variance*, the square root of which is interpreted as an “error bar”. Such uncertainty estimates have numerous applications, one example being active learning, where one explores input locations where uncertainties over output values are high.

The posterior variance is, by definition, the posterior expected square difference between the posterior mean $\bar{m}(x)$ given by (16) and the output $f(x)$ of posterior GP sample $f \sim \mathcal{GP}(\bar{m}, \bar{k})$, that is,

$$\bar{k}(x, x) = \mathbb{E}_{f \sim \mathcal{GP}(\bar{m}, \bar{k})} [(f(x) - \bar{m}(x))^2]. \quad (33)$$

In other words, $\bar{k}(x, x)$ is interpreted as the *average case error* at a location x from the Bayesian viewpoint. The purpose of this subsection is show that there exists a kernel/frequentist interpretation of $\bar{k}(x, x)$ as a certain *worst case error*. To the best of our knowledge, this fact has not been known in the literature.

For simplicity, we focus here on regression with a zero-mean GP prior $f \sim \mathcal{GP}(0, k)$. We use the following notation. Let $w^\sigma : \mathcal{X} \rightarrow \mathbb{R}^n$ be a vector-valued function defined by

$$w^\sigma(x) := (k_{XX} + \sigma^2 I_n)^{-1} k_{Xx} \in \mathbb{R}^n, \quad x \in \mathcal{X}, \quad (34)$$

so that the posterior mean function (16) can be written as a projection of $Y = (y_1, \dots, y_n)^\top$ onto $w^\sigma(x)$:

$$\bar{m}(x) = \sum_{i=1}^n w_i^\sigma(x) y_i = Y^\top w^\sigma(x). \quad (35)$$

Moreover, define a new kernel $k^\sigma : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ by

$$k^\sigma(x, y) := k(x, y) + \sigma^2 \delta(x, y), \quad x, y \in \mathcal{X}, \quad (36)$$

where $\delta : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ is the kernel defined in (31). Then, by definition, the kernel matrix of k^σ for $X = (x_1, \dots, x_n)$ is given by $k_{XX}^\sigma = k_{XX} + \sigma^2 I_n$. Note that (36) is understood as the covariance kernel of the contaminated observation process (32), where $\mathbf{f} := \mathbf{f} \sim \mathcal{GP}(0, k)$. Let \mathcal{H}_{k^σ} be the RKHS of k^σ .

Given these preliminaries, we now present our result on the worst case error interpretation of the posterior variance (33).

Proposition 3.8 *Let \bar{k} be the posterior covariance function (17) with noise variance σ^2 . Then, for any $x \in \mathcal{X}$ with $x \neq x_i$, $i = 1, \dots, n$, we have*

$$\sqrt{\bar{k}(x, x) + \sigma^2} = \sup_{g \in \mathcal{H}_{k^\sigma} : \|g\|_{\mathcal{H}_{k^\sigma}} \leq 1} \left(g(x) - \sum_{i=1}^n w_i^\sigma(x) g(x_i) \right). \quad (37)$$

To prove the above proposition, we need the following lemma, which is useful in general. The proof can be found in Appendix A.1

Lemma 3.9 *Let k be a kernel on \mathcal{X} and \mathcal{H}_k be its RKHS. Then for any $m \in \mathbb{N}$, $x_1, \dots, x_m \in \mathcal{X}$ and $c_1, \dots, c_m \in \mathbb{R}$, we have*

$$\left\| \sum_{i=1}^m c_i k(\cdot, x_i) \right\|_{\mathcal{H}_k} = \sup_{f \in \mathcal{H}_k : \|f\|_{\mathcal{H}_k} \leq 1} \sum_{i=1}^m c_i f(x_i). \quad (38)$$

Lemma 3.9 shows that the RKHS norm of a linear combination of feature vectors can be written as a supremum over functions in the unit ball of the RKHS. Based on this result, Proposition 3.8 can be proven as follows.

Proof By Lemma 3.9, we have

$$\left\| k^\sigma(\cdot, x) - \sum_{i=1}^n w_i^\sigma(x) k^\sigma(\cdot, x_i) \right\|_{\mathcal{H}_{k^\sigma}} = \sup_{\substack{g \in \mathcal{H}_{k^\sigma} \\ \|g\|_{\mathcal{H}_{k^\sigma}} \leq 1}} \left(g(x) - \sum_{i=1}^n w_i^\sigma(x) g(x_i) \right). \quad (39)$$

The left side of this equality can be expanded as

$$\begin{aligned}
& \left\| k^\sigma(\cdot, x) - \sum_{i=1}^n w_i^\sigma(x) k^\sigma(\cdot, x_i) \right\|_{\mathcal{H}_{k^\sigma}}^2 \\
&= k^\sigma(x, x) - 2 \sum_{i=1}^n w_i^\sigma(x) k^\sigma(x, x_i) + \sum_{i,j=1}^n w_i^\sigma(x) w_j^\sigma(x) k^\sigma(x_i, x_j) \\
&= k(x, x) + \sigma^2 - 2 \sum_{i=1}^n w_i^\sigma(x) k(x, x_i) + w^\sigma(x)^\top k_{XX}^\sigma w^\sigma(x) \\
&= k(x, x) + \sigma^2 - 2w^\sigma(x)^\top k_{Xx} + k_{xX} (k_{XX} + \sigma^2 I_n)^{-1} (k_{XX} + \sigma^2 I_n) (k_{XX} + \sigma^2 I_n)^{-1} k_{Xx} \\
&= k(x, x) + \sigma^2 - 2k_{xX} (k_{XX} + \sigma^2 I_n)^{-1} k_{Xx} + k_{xX} (k_{XX} + \sigma^2 I_n)^{-1} k_{Xx} \\
&= k(x, x) + \sigma^2 - k_{xX} (k_{XX} + \sigma^2 I_n)^{-1} k_{Xx} \\
&= \bar{k}(x, x) + \sigma^2,
\end{aligned}$$

where we used the assumption $x \neq x_i$ for $i = 1, \dots, n$ in the second equality. The assertion follows from this last expression and (39). \blacksquare

Remark 3.10 To understand Proposition 3.8, we need to understand the structure of the RKHS \mathcal{H}_{k^σ} . First, the RKHS of $\sigma^2 \delta$, denoted by $\mathcal{H}_{\sigma^2 \delta}$, is characterized as

$$\mathcal{H}_{\sigma^2 \delta} = \left\{ h = \sum_{x \in \mathcal{X}} c_x \sigma^2 \delta(\cdot, x) : \|h\|_{\mathcal{H}_{\sigma^2 \delta}}^2 = \sigma^4 \sum_{x \in \mathcal{X}} c_x^2 < \infty, \quad c_x \in \mathbb{R}, \quad x \in \mathcal{X} \right\}, \quad (40)$$

where the summation $\sum_{x \in \mathcal{X}}$ can be uncountable. It is known that $\mathcal{H}_{\sigma^2 \delta}$ is not separable; see e.g. Steinwart and Scovel [2012, Example 3.9]. Note that the kernel $\sigma^2 \delta(x, y)$ is given as a multiplication of σ^2 to the kernel $\delta(x, y)$, so $\mathcal{H}_{\sigma^2 \delta}$ is norm-equivalent to the RKHS \mathcal{H}_δ of δ . Moreover, from (40), it is easy to see that for any $h \in \mathcal{H}_\delta$, we have $\|h\|_{\mathcal{H}_{\sigma^2 \delta}} = \|h\|_{\mathcal{H}_\delta}$ for all $\sigma > 0$.

Since k^σ is defined as the sum of two kernels k and $\sigma^2 \delta$, the RKHS \mathcal{H}_{k^σ} consists of functions that can be written as a sum of functions from \mathcal{H}_k and $\mathcal{H}_{\sigma^2 \delta}$ [Aronszajn, 1950, Section 6]:

$$\mathcal{H}_{k^\sigma} = \{g = f + h : f \in \mathcal{H}_k, h \in \mathcal{H}_{\sigma^2 \delta}\} = \{g = f + h : f \in \mathcal{H}_k, h \in \mathcal{H}_\delta\}, \quad (41)$$

where the corresponding RKHS norm is given by

$$\|g\|_{\mathcal{H}_{k^\sigma}} = \inf_{\substack{f \in \mathcal{H}_k, h \in \mathcal{H}_{\sigma^2 \delta} \\ g=f+h}} \|f\|_{\mathcal{H}_k} + \|h\|_{\mathcal{H}_{\sigma^2 \delta}} = \inf_{\substack{f \in \mathcal{H}_k, h \in \mathcal{H}_\delta \\ g=f+h}} \|f\|_{\mathcal{H}_k} + \|h\|_{\mathcal{H}_\delta}.$$

Remark 3.11 In (37), the quantity $\sum_{i=1}^n w_i^\sigma(x) g(x_i)$ for a fixed $g \in \mathcal{H}_{k^\sigma}$ with $\|g\|_{\mathcal{H}_{k^\sigma}} \leq 1$ is the posterior mean function given training data $(x_i, g(x_i))_{i=1}^n$. Note that by (41), g can be written as $g = f + h$ for some $f \in \mathcal{H}_k$ and $h \in \mathcal{H}_{\sigma^2 \delta}$ such that $\|f\|_{\mathcal{H}_k} + \|h\|_{\mathcal{H}_{\sigma^2 \delta}} \leq 1$. Therefore each training output $g(x_i)$ can be written as $g(x_i) = f(x_i) + h(x_i)$, where $h(x_i)$ is understood as “independent noise.” The supremum in (37) is thus the worst case error between the “true” value $g(x) = f(x) + h(x)$ and the estimate $\sum_{i=1}^n w_i^\sigma(x) g(x_i) = \sum_{i=1}^n w_i^\sigma(x) (f(x_i) + h(x_i))$.

Note that since $h(x)$ is “independent noise,” it is not possible to estimate this value from the training data, and thus a certain amount of error is unavoidable. This may intuitively explain why the additional σ^2 term appears in the left side of (37), which shows that the worst case error is more pessimistic than the average case error, in the presence of noise.

Remark 3.12 From the proof of Lemma 3.9, it is easy to see that the function $g \in \mathcal{H}_{k^\sigma}$ that attains the supremum in (37) is

$$\begin{aligned} g &= C \left(k^\sigma(\cdot, x) - \sum_{i=1}^n w_i^\sigma(x) k^\sigma(\cdot, x_i) \right) \\ &= C \left(k(\cdot, x) - \sum_{i=1}^n w_i^\sigma(x) k(\cdot, x_i) \right) + C\sigma^2 \left(\delta(\cdot, x) - \sum_{i=1}^n w_i^\sigma(x) \delta(\cdot, x_i) \right). \end{aligned}$$

where $C := \|k^\sigma(\cdot, x) - \sum_{i=1}^n w_i^\sigma(x) k^\sigma(\cdot, x_i)\|_{\mathcal{H}_{k^\sigma}}^{-1}$ is a normalizing constant that ensures that $\|f\|_{\mathcal{H}_{k^\sigma}} = 1$. Thus, the worst adversarial function can be written as $g = f + h$, where

$$f := C \left(k(\cdot, x) - \sum_{i=1}^n w_i^\sigma(x) k(\cdot, x_i) \right) \in \mathcal{H}_k \quad \text{and} \quad h := C\sigma^2 \left(\delta(\cdot, x) - \sum_{i=1}^n w_i^\sigma(x) \delta(\cdot, x_i) \right) \in \mathcal{H}_{\sigma^2 \delta}.$$

This implies that, as the noise variance σ^2 increases, the relative contribution of the noise term h to the adversarial function g increases; this makes it more difficult to fit the adversarial function, and thus the worst case error becomes more pessimistic than the average case error, as shown in (37).

Noise-free case. We consider the important special case of noise-free observations, that is the case where $\sigma^2 = 0$, to further illustrate Proposition 3.8. In this case, the posterior standard deviation $\sqrt{\bar{k}(x, x)}$, or (the square-root of) the average case error, is identical to the worst case error; see Fig. 2 for illustration. The following result, which does not require $x \neq x_i$ for $i = 1, \dots, n$ as opposed to Proposition 3.8, can be proven in a similar way to that of Proposition 3.8.

Proposition 3.10 *Assume that $\sigma^2 = 0$, and that the kernel matrix k_{XX} is invertible. Then we have*

$$\sqrt{\bar{k}(x, x)} = \sup_{\substack{f \in \mathcal{H}_k \\ \|f\|_{\mathcal{H}_k} \leq 1}} \left(\sum_{i=1}^n w_i(x) f(x_i) - f(x) \right), \quad x \in \mathcal{X}, \quad (42)$$

where $(w_1(x), \dots, w_n(x))^\top = k_{XX}^{-1} k_{Xx} \in \mathbb{R}^n$ and \bar{k} is given by (17) with $\sigma^2 = 0$.

By applying the Cauchy-Schwartz inequality to (42), we have the following corollary. It shows that the posterior variance $\bar{k}(x, x)$ provides an upper-bound on the error of kernel-based interpolation for a fixed target function.

Corollary 3.11 *Assume that $\sigma^2 = 0$, and that the kernel matrix k_{XX} is invertible. Then for all $f \in \mathcal{H}_k$, we have*

$$(\bar{m}(x) - f(x))^2 \leq \|f\|_{\mathcal{H}_k}^2 \bar{k}(x, x), \quad x \in \mathcal{X}.$$

where $\bar{m}(x) = \sum_{i=1}^n w_i(x) f(x_i)$.

3.5 A Weight Vector Viewpoint of Regularization and the Additive Noise Assumption

Based on the weight vector (34) and the worst case error in the right side of (42), we provide another interpretation of the equivalence between regularization and the additive noise assumption in regression. This is given by the following result.

Proposition 3.12 *Let $X = (x_1, \dots, x_n) \in \mathcal{X}^n$ be fixed, and let $Y = (y_1, \dots, y_n)^\top \in \mathbb{R}^n$ with $y_i = f(x_i) + \xi_i$, $i = 1, \dots, n$, where $f : \mathcal{X} \rightarrow \mathbb{R}$ is a fixed function and ξ_1, \dots, ξ_n are random variables such that $\mathbb{E}[\xi_i] = 0$ and $\mathbb{E}[\xi_i \xi_j] = \sigma^2 \delta_{ij}$ for $\sigma > 0$. Let $w^\sigma : \mathcal{X} \rightarrow \mathbb{R}^n$ be the vector-valued function as defined in (34) with σ , X and kernel $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$. Then we have*

$$w^\sigma(x) = \arg \min_{w \in \mathbb{R}^n} \left[\sup_{\|f\|_{\mathcal{H}_k} \leq 1} (f(x) - f_X^\top w) \right]^2 + \sigma^2 \|w\|^2 \quad (43)$$

$$= \arg \min_{w \in \mathbb{R}^n} \left[\sup_{\|f\|_{\mathcal{H}_k} \leq 1} (f(x) - f_X^\top w) \right]^2 + \text{var}_{\xi_1, \dots, \xi_n} [Y^\top w], \quad x \in \mathcal{X}, \quad (44)$$

where $f_X := (f(x_1), \dots, f(x_n))^\top \in \mathbb{R}^n$.

Proof First, by the reproducing property it is easy to show that

$$w^\sigma(x) = \arg \min_{w \in \mathbb{R}^n} \left\| k(\cdot, x) - \sum_{i=1}^n w_i k(\cdot, x_i) \right\|_{\mathcal{H}_k}^2 + \sigma^2 \|w\|^2, \quad x \in \mathcal{X}. \quad (45)$$

For the first term in the right side, we have by Lemma 3.9,

$$\left\| k(\cdot, x) - \sum_{i=1}^n w_i k(\cdot, x_i) \right\|_{\mathcal{H}_k} = \sup_{\|f\|_{\mathcal{H}_k} \leq 1} (f(x) - f_X^\top w).$$

On the other hand, for the second term in the right side of (45), we have

$$\text{var}_{\xi_1, \dots, \xi_n} [Y^\top w] = \mathbb{E}_{\xi_1, \dots, \xi_n} [(\xi^\top w)^2] = \mathbb{E}_{\xi_1, \dots, \xi_n} [w^\top \xi \xi^\top w] = \sigma^2 \|w\|^2, \quad (46)$$

where $\xi := (\xi_1, \dots, \xi_n)^\top \in \mathbb{R}^n$. The assertion follows by inserting these identities in (45). \blacksquare

Remark 3.13 Note that in (44) and (46), the noise variables ξ_1, \dots, ξ_n are independent of the function values $f(x_1), \dots, f(x_n)$, because of our assumption in Proposition 3.12; note also that the function f is fixed, and is not assumed to be a Gaussian process.

Remark 3.14 To discuss Proposition 3.12, let us fix a weight vector $w \in \mathbb{R}^n$. Then the first term in the right side of (43) is the worst case error in the noise-free setting, since $f_X^\top w = \sum_{i=1}^n w_i f(x_i)$ can be considered as an estimator of $f(x)$ based on noise-free observations $f(x_1), \dots, f(x_n)$, where f is taken from the unit ball in the RKHS; recall also (35). On the other hand, the second term in (43) is a regularizer that makes the squared Euclidian norm of the weight vector w not too large. Importantly, it shows that the noise variance σ^2 serves as a regularization constant.

Regarding the second term in (44), $Y^\top w$ is an estimator of $f(x)$, where f is the (fixed) latent regression function; see again (35). Thus $\text{var}_{\xi_1, \dots, \xi_n}[Y^\top w]$ is the variance of the regression estimator, which is equal to the regularization term $\sigma^2 \|w\|^2$ in (43). Therefore, (44) shows that the weight vector (34) is obtained so as to minimize the sum of the noise-free worst case error and the variance of the regression estimator based on noisy observations.

4 Hypothesis Spaces: Do Gaussian Process Draws Lie in an RKHS?

In discussions about the similarity between GPs and kernel methods, it is often pointed out that the hypothesis space of Gaussian processes is not equal to that of kernel ridge regression (i.e., the corresponding RKHS). For instance, Neal [1998, Section 7] discussed this topic, arguing why GP models had not been widely used at the time of his writing:

I speculate that a more fundamental reason for the neglect of Gaussian process models is a widespread preference for simple models, resulting from a confusion between prior beliefs regarding the true function being modeled and expectations regarding the properties of the best predictor for this function (the posterior mean, under squared error loss). These need not be at all similar. For example, our beliefs about the true function might sometimes be captured by an Ornstein-Uhlenbeck process, a Gaussian process with covariance function $\exp(-|x^{(i)} - x^{(j)}|)$. Realizations from this process are nowhere differentiable, but the predictive mean function will consist of pieces that are sums of exponentials, as can be seen from equation (3).

As explained in Section 3.3, the posterior mean function of GP-regression (which is the “predictive mean function” in the above quotation) lies in the RKHS of the GP covariance kernel. On the other hand, if one considers a sample path of the GP prior (the Ornstein-Uhlenbeck process in the quotation, which is the GP of the Matérn kernel with $\alpha = 1/2$ and $d = 1$; see also Example 2.5), it is almost surely less smooth than functions in the RKHS (which is norm-equivalent to the first-order Sobolev space; see Example 2.6.) and hence does belong to that RKHS almost surely. This is the difference Neal [1998] mentioned.

The purpose of this section is to explain why the above mentioned difference exists, by reviewing sample path properties of GPs and how they are related to RKHSs. To this end, in Section 4.1 we first look at characterizations of GPs and RKHSs by orthonormal expansions, namely the *Karhunen-Loève expansion* for GPs and *Mercer representation* for RKHSs. We then review Driscoll’s theorem [Driscoll, 1973, Lukić and Beder, 2001], which provides a necessary and sufficient condition for a GP sample path to lie in a *given* RKHS (which can be different from the RKHS associated with the GP covariance kernel) in Section 4.2. Using this result, we show in Section 4.3 that GP sample spaces can be constructed as *powers of RKHSs* defined from the Mercer representation [Steinwart and Scovel, 2012]; this recovers a special case of recent generic results by Steinwart [2017] on sample path properties. We conclude this section by using these results to derive GP sample path properties for square-exponential kernels and Matérn kernels in Section 4.4, the latter providing a theoretical explanation of the above difference mentioned by Neal [1998].

The main message of this section may be summarized as follows: While GP sample paths fall outside of the RKHS of the GP covariance kernel almost surely, they actually lie on certain

RKHSs defined as powers of that RKHS; therefore GPs and RKHSs are still deeply connected in terms of the induced hypothesis spaces, and the difference such as the one mentioned by Neal [1998] should not warrant strong conceptual separation between the two frameworks. We will also use the sample path properties in this section to discuss the equivalence between convergence properties of GP and kernel ridge regression in Section 5.

4.1 Characterizations via Orthonormal Expansions

To gain intuition and understanding on the structure of RKHS and GP, we review here their expressions via orthonormal functions, that is, *Karhunen-Loève expansion* for GPs and *Mercer representation* for RKHSs. These expressions are given in terms of the eigenvalues and eigenfunctions of a kernel integral operator defined below. For simplicity, we assume here that \mathcal{X} is a compact metric space (e.g., a bounded and closed subset of \mathbb{R}^d), and k is a continuous kernel on \mathcal{X} .

4.1.1 Mercer's Theorem

Let ν be a finite Borel measure on \mathcal{X} with \mathcal{X} being its support (e.g., the Lebesgue measure on $\mathcal{X} \subset \mathbb{R}^d$). Let $L_2(\nu)$ be the Hilbert space of square-integrable functions⁵ with respect to ν , as defined in (1) with $p = 2$. Define an operator $T_k : L_2(\nu) \rightarrow L_2(\nu)$ as the integral operator with the kernel k and the measure ν :

$$T_k f := \int k(\cdot, x) f(x) d\nu(x), \quad f \in L_2(\nu). \quad (47)$$

If the kernel k is defined on $\mathcal{X} \subset \mathbb{R}^d$ and shift-invariant (that is, it can be written in the form $k(x, y) = \phi(x - y)$ for some positive definite function ϕ), then this operator is a convolution of k and a function f .⁶ Therefore the output function $T_k f$ can be seen as a smoothed version of f , if k is smooth.

Since T_k is compact, positive and self-adjoint, the spectral theorem (see e.g. Steinwart and Christmann 2008, Theorem A.5.13) guarantees an eigen-decomposition of T_k in the form

$$T_k f = \sum_{i \in I} \lambda_i \langle \phi_i, f \rangle_{L_2(\nu)} \phi_i, \quad (48)$$

where the convergence is in $L_2(\nu)$. Here $I \subset \mathbb{N}$ is a set of indices (e.g., $I = \mathbb{N}$ when the RKHS is infinite dimensional, and $I = \{1, \dots, K\}$ with $K \in \mathbb{N}$ when the RKHS is K -dimensional), and $(\phi_i, \lambda_i)_{i \in I} \subset L_2(\nu) \times (0, \infty)$ are (countable) eigenfunctions and the associated eigenvalues of T_k such that $\lambda_1 \geq \lambda_2 \geq \dots > 0$:

$$T_k \phi_i = \lambda_i \phi_i, \quad i \in I.$$

The eigenfunctions $(\phi_i)_{i \in \mathbb{N}}$ form an orthonormal system in $L_2(\nu)$, i.e., $\langle \phi_i, \phi_j \rangle_{L_2(\nu)} = \delta_{ij}$, where $\delta_{ij} = 1$ if $i = j$ and $\delta_{ij} = 0$ otherwise.

Mercer's theorem, which is named after Mercer [1909], states the kernel k can be expressed in terms of the eigensystem $(\phi_i, \lambda_i)_{i \in I}$ in (48). This expression of the kernel provides useful

⁵Strictly, here each $f \in L_2(\nu)$ represents the class of functions that are equivalent ν -almost everywhere.

⁶Or more precisely, a convolution of the positive definite function ϕ and a measure $d\eta := f d\nu$ that has f as a Radon-Nikodym derivative w.r.t. ν

ways of constructing GPs and RKHSs, as described shortly. The following form of Mercer's theorem is due to Steinwart and Christmann [2008, Theorem 4.49], while we note that Mercer's theorem holds under weaker assumptions than those considered here [Steinwart and Scovel, 2012, Section 3].

Theorem 4.1 (Mercer's theorem) *Let \mathcal{X} be a compact metric space, $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ be a continuous kernel, ν be a finite Borel measure whose support is \mathcal{X} , and $(\phi_i, \lambda_i)_{i \in I}$ be as in (48). Then we have*

$$k(x, x') = \sum_{i \in I} \lambda_i \phi_i(x) \phi_i(x'), \quad x, x' \in \mathcal{X}, \quad (49)$$

where the convergence is absolute and uniform over $x, x' \in \mathcal{X}$.

Remark 4.1 The expansion in (49) depends on the measure ν , since $(\phi_i, \lambda_i)_{i \in I}$ is an eigensystem of the integral operator (47), which is defined with ν . However, the kernel k in the left side is unique, irrespective of the choice of ν . In other words, a different choice of ν results in a different eigensystem $(\phi_i, \lambda_i)_{i \in I}$, and thus results in a different basis expression of the same kernel k .

Remark 4.2 In Theorem 4.1, the assumption that ν has \mathcal{X} as its support is important, since otherwise the equality (49) may not hold for some $x \in \mathcal{X}$. For instance, assume that there is an open set $N \subset \mathcal{X}$ such that $\nu(N) = 0$. Then the integral operator (47) does not take into account the values of a function f on N , and therefore the eigenfunctions ϕ_i are only uniquely defined on $\mathcal{X} \setminus N$, in which case, the equality in (49) holds only on $\mathcal{X} \setminus N$. We refer to Steinwart and Scovel [2012, Corollaries 3.2 and 3.5] for precise statements of Mercer's theorem in such a case.

4.1.2 Mercer Representation of RKHSs

The eigensystem of the integral operator (47) provides a series representation of the RKHS, which is called the Mercer representation [Steinwart and Christmann, 2008, Theorem 4.51].

Theorem 4.2 (Mercer Representation) *Let \mathcal{X} be a compact metric space, $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ be a continuous kernel, ν be a finite Borel measure whose support is \mathcal{X} , and $(\phi_i, \lambda_i)_{i \in I}$ be as in (48). Then the RKHS \mathcal{H}_k of k is given by*

$$\mathcal{H}_k = \left\{ f := \sum_{i \in I} \alpha_i \lambda_i^{1/2} \phi_i : \|f\|_{\mathcal{H}_k}^2 := \sum_{i \in I} \alpha_i^2 < \infty \right\}, \quad (50)$$

and the inner-product is given by

$$\langle f, g \rangle_{\mathcal{H}_k} = \sum_{i \in I} \alpha_i \beta_i \quad \text{for} \quad f := \sum_{i \in I} \alpha_i \lambda_i^{1/2} \phi_i \in \mathcal{H}_k, \quad g := \sum_{i \in I} \beta_i \lambda_i^{1/2} \phi_i \in \mathcal{H}_k.$$

In other words, $(\lambda_i^{1/2} \phi)_{i \in I}$ forms an orthonormal basis of \mathcal{H}_k .

Remark 4.3 As mentioned in Remark 4.1 for the expansion of a kernel (49), the Mercer representation in (50) depends on the measure ν , since $(\phi_i, \lambda_i)_{i \in I}$ depends on ν . Under the assumptions in Theorem 4.2, however, a different choice of ν , which results in a different eigensystem $(\phi_i, \lambda_i)_{i \in I}$, results in the same RKHS \mathcal{H}_k . Note also here that, as mentioned in Remark 4.2, the assumption that ν has its support on \mathcal{X} is crucial.

4.1.3 Karhunen-Loève Expansion of Gaussian Processes

Corresponding to the Mercer representation of RKHSs, there exists a series representation of Gaussian processes known as the *Karhunen-Loève (KL) expansion*. The KL-expansion is based on the eigensystem of the integral operator in (47), as for the Mercer representation. This is a consequence of the canonical isometric isomorphism between an RKHS and the corresponding *Gaussian Hilbert space* [Janson, 1997]. The following result, which is well known in the literature, follows from Steinwart [2017, Lemmas 3.3 and 3.7]; see also e.g., Adler [1990, Sections 3.2 and 3.3] and Berlinet and Thomas-Agnan [2004, Section 2.3].

Theorem 4.3 (Karhunen-Loève Expansion) *Let \mathcal{X} be a compact metric space, $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ be a continuous kernel, ν be a finite Borel measure whose support is \mathcal{X} , and $(\phi_i, \lambda_i)_{i \in I}$ be as in (48). For a Gaussian process $f \sim \mathcal{GP}(0, k)$, define*

$$\mathbf{z}_i := \lambda_i^{-1/2} \int f(x) \phi_i(x) d\nu(x), \quad i \in I. \quad (51)$$

Then the following are true:

1. We have

$$\mathbf{z}_i \sim \mathcal{N}(0, 1) \quad \text{and} \quad \mathbb{E}[\mathbf{z}_i \mathbf{z}_j] = \delta_{ij}, \quad i, j \in I. \quad (52)$$

2. For all $x \in \mathcal{X}$ and for all finite $J \subset I$, we have

$$\mathbb{E} \left[\left(f(x) - \sum_{i \in J} \mathbf{z}_i \lambda_i^{1/2} \phi_i(x) \right)^2 \right] = k(x, x) - \sum_{j \in J} \lambda_j e_j^2(x). \quad (53)$$

3. If $I = \mathbb{N}$, we have

$$\lim_{n \rightarrow \infty} \mathbb{E} \left[\left(f(x) - \sum_{i=1}^n \mathbf{z}_i \lambda_i^{1/2} \phi_i(x) \right)^2 \right] = 0, \quad x \in \mathcal{X}, \quad (54)$$

where the convergence is uniform in $x \in \mathcal{X}$.

Remark 4.4 Informally, Theorem 4.3 shows that the Gaussian process $f \sim \mathcal{GP}(0, k)$ can be expressed using ONB $(\lambda_i^{1/2} \phi_i)_{i \in I}$ of \mathcal{H}_k and standard Gaussian random variables $\mathbf{z}_i \sim \mathcal{N}(0, 1)$ as

$$f(x) = \sum_{i \in I} \mathbf{z}_i \lambda_i^{1/2} \phi_i(x), \quad x \in \mathcal{X} \quad (55)$$

where the convergence is in the mean square sense and uniform over $x \in \mathcal{X}$, as shown in (54). Note that (54) is an immediate consequence of (53) and Mercer's theorem (Theorem 4.1). Steinwart [2017, Theorem 3.5] shows that, under the same conditions, the convergence in (55) also holds in $L_2(\nu)$. The expression (55) is what is often called the KL expansion.

Remark 4.5 (52) shows that $(\mathbf{z}_i)_{i \in I}$ as defined in (51) are standard normal, and are independent to each other. Note that $(\mathbf{z}_i)_{i \in I}$ are dependent to the given $f \sim \mathcal{GP}(0, k)$, as can be seen from (51); otherwise (53) and (54) do not hold. On the other hand, *given* i.i.d. standard normal random variables $\mathbf{z}_1, \dots, \mathbf{z}_n \stackrel{i.i.d.}{\sim} \mathcal{N}(0, 1)$ (i.e., independent to a *specific* realization $f \sim \mathcal{GP}(0, k)$), one can construct a finite dimensional Gaussian process in the form $\sum_{i=1}^n \mathbf{z}_i \lambda_i^{1/2} \phi_i$ that is approximately distributed as $\mathcal{GP}(0, k)$; this is often called a *truncated* KL expansion.

4.2 Sample Path Properties and the Zero-One Law

We review Driscoll's theorem, which provides a necessary and sufficient condition for a Gaussian process $f \sim \mathcal{GP}(0, k)$ to belong to an RKHS \mathcal{H}_r with kernel r with probability 1 or 0 [Driscoll, 1973, Theorem 3]. Here the kernels r and k can be different in general, but are defined on the same space \mathcal{X} . Since these probabilities (i.e., 1 or 0) are the only options, the theorem is called Driscoll's zero-one law. (In other words, a statement like " $f \sim \mathcal{GP}(0, k)$ belongs to \mathcal{H}_r with probability 0.3" is false.) We review in particular a generalization of Driscoll's theorem by Lukić and Beder [2001, Theorem 7.4], which holds under weaker assumptions than the original theorem by Driscoll [1973]. Our presentation below also uses some facts pointed out by Steinwart [2017]. To state the result of Lukić and Beder [2001], we need to introduce the notion of the *dominance operator*, whose existence is shown by Lukić and Beder [2001, Theorem 1.1].

Theorem 4.4 (Dominance operator) *Let k and r be positive definite kernels on a set \mathcal{X} , and let \mathcal{H}_k and \mathcal{H}_r be their respective RKHSs. Assume $\mathcal{H}_k \subset \mathcal{H}_r$, and let $I_{kr} : \mathcal{H}_k \rightarrow \mathcal{H}_r$ be the natural inclusion operator, i.e., $I_{kr}g := g$ for $g \in \mathcal{H}_k$. Then I_{kr} is continuous. Moreover, there exists a unique linear operator $L : \mathcal{H}_r \rightarrow \mathcal{H}_k$ such that*

$$\langle f, g \rangle_{\mathcal{H}_r} = \langle Lf, g \rangle_{\mathcal{H}_k}, \quad \forall f \in \mathcal{H}_r, \forall g \in \mathcal{H}_k. \quad (56)$$

In particular, we have

$$Lr(\cdot, x) = k(\cdot, x), \quad \forall x \in \mathcal{X}.$$

Furthermore, $I_{kr}L : \mathcal{H}_r \rightarrow \mathcal{H}_r$ is bounded, positive and symmetric.

The key concept in Driscoll's theorem is the *nuclear dominance*, which is defined in the following way. As explained shortly, the nuclear dominance serves as a necessary and sufficient condition in Driscoll's zero-one law.

Definition 4.5 (Nuclear dominance) *Under the same notation as in Theorem 4.4, assume $\mathcal{H}_k \subset \mathcal{H}_r$. Then r is said to dominate k , and the operator L in Theorem 4.4 is called the dominance operator of \mathcal{H}_r over \mathcal{H}_k . Moreover, the dominance is called nuclear, in which case it is written as $r \gg k$, if $I_{kr}L : \mathcal{H}_r \rightarrow \mathcal{H}_r$ is nuclear (or of trace class), i.e.,*

$$\text{Tr}(I_{kr}L) = \sum_{i \in I} \langle I_{kr}L\psi_i, \psi_i \rangle_{\mathcal{H}_r} < \infty,$$

where $(\psi_i)_{i \in I} \subset \mathcal{H}_r$ is an ONB of \mathcal{H}_r .

Before stating Driscoll's theorem, we mention that the dominance operator in Theorem 4.4 can be written in terms of the inclusion operator I_{kr} . That is, Steinwart [2017, Section 2] pointed out that the dominance operator L is identical to I_{kr}^* , the adjoint operator of I_{kr} , as summarized in the following lemma.

Lemma 4.6 *Under the same notation as in Theorem 4.4, assume $\mathcal{H}_k \subset \mathcal{H}_r$. Let L be the dominance operator as given in Theorem 4.4. Then we have $L = I_{kr}^*$.*

Proof Let I_{kr}^* be the adjoint of I_{kr} . Then we have

$$\langle g, f \rangle_{\mathcal{H}_r} = \langle I_{kr}g, f \rangle_{\mathcal{H}_r} = \langle g, I_{kr}^*f \rangle_{\mathcal{H}_r}, \quad \forall f \in \mathcal{H}_r, \forall g \in \mathcal{H}_k,$$

which is the property (56) of the dominance operator. Since the dominance operator is unique by Theorem 4.4, we have $L = I_{kr}^*$. ■

The following result shows that the nuclear dominance is equivalent to the inclusion operator I_{kr} being Hilbert-Schmidt. The result is essentially given by the proof of Steinwart [2017, Lemma 7.4, equivalence of (i) and (iii)].

Lemma 4.7 *Under the same notation as in Theorem 4.4, assume $\mathcal{H}_k \subset \mathcal{H}_r$. Then the following statements are equivalent:*

1. *The nuclear dominance holds: $r \gg k$, i.e., $I_{kr}L : \mathcal{H}_r \rightarrow \mathcal{H}_r$ is nuclear.*
2. *The inclusion operator $I_{kr} : \mathcal{H}_k \rightarrow \mathcal{H}_r$ is Hilbert-Schmidt.*

Proof From Lemma 4.6, we have

$$\mathrm{Tr}(I_{kr}L) = \mathrm{Tr}(I_{kr}I_{kr}^*) := \sum_{i \in I} \langle I_{kr}I_{kr}^*\psi_i, \psi_i \rangle_{\mathcal{H}_r} = \sum_{i \in I} \langle I_{kr}^*\psi_i, I_{kr}^*\psi_i \rangle_{\mathcal{H}_r} =: \|I_{kr}^*\|_{\mathrm{HS}}^2,$$

where $\|\cdot\|_{\mathrm{HS}}$ denotes the Hilbert-Schmidt norm and $\mathrm{Tr}(\cdot)$ the trace, and $(\psi_i)_{i \in I} \subset \mathcal{H}_r$ is an ONB of \mathcal{H}_r . Since we have $\|I_{kr}^*\|_{\mathrm{HS}} = \|I_{kr}\|_{\mathrm{HS}}$ (see e.g., Steinwart and Christmann 2008, p.506), the assertion immediately follows. ■

Using Lemma 4.7, Theorem 7.4 of Lukić and Beder [2001], which is a generalization of the zero-one law of Driscoll [1973, Theorem 3], can be stated as Theorem 4.9 below. To state it, we need to introduce a definition of a stochastic process being a *version* of a GP [Brémaud, 2014, Definition 3.1.9].

Definition 4.8 (A version of a GP) *Let $f \sim \mathcal{GP}(m, k)$ be a Gaussian process with mean function $m : \mathcal{X} \rightarrow \mathbb{R}$ and covariance kernel $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$, where \mathcal{X} is a nonempty set. Then a stochastic process \tilde{f} on \mathcal{X} is called a version of f , if $f(x) = \tilde{f}(x)$ holds with probability 1 for all $x \in \mathcal{X}$.*

Theorem 4.9 (A generalized Driscoll's theorem) *Let k and r be positive definite kernels on a set \mathcal{X} , and let \mathcal{H}_k and \mathcal{H}_r be their respective RKHSs. Assume $\mathcal{H}_k \subset \mathcal{H}_r$, and let $I_{kr} : \mathcal{H}_k \rightarrow \mathcal{H}_r$ be the natural inclusion operator. Let $f \sim \mathcal{GP}(m, k)$ be a Gaussian process such that $m \in \mathcal{H}_r$. Then the following statements are true.*

1. *If I_{kr} is Hilbert-Schmidt, then there is a version \tilde{f} of f such that $\tilde{f} \in \mathcal{H}_r$ holds with probability 1.*
2. *If I_{kr} is not Hilbert-Schmidt, then $f \in \mathcal{H}_r$ holds with probability 0.*

Remark 4.6 In Driscoll [1973, Theorem 3], it is assumed that \mathcal{X} is a separable metric space, k is a continuous kernel on \mathcal{X} and $f \sim \mathcal{GP}(m, k)$ is almost surely continuous. Under this assumption, Driscoll [1973, Theorem 3] showed that a condition equivalent to the nuclear dominance condition [Lukić and Beder, 2001, Proposition 4.5] implies that the *given* Gaussian process f belongs to \mathcal{H}_r with probability 1. That is, in this case one does not need to consider a version \tilde{f} of it.

Remark 4.7 In Lukić and Beder [2001, Theorem 5.1], it is shown that the nuclear dominance condition (which is equivalent to I_{kr} being Hilbert-Schmidt) implies that *any* second-order process f with covariance kernel k (i.e., f does not necessarily be Gaussian) belongs to \mathcal{H}_r with probability 1.

Remark 4.8 One way to check whether I_{kr} is Hilbert-Schmidt is given by Gonzalez-Barrios and Dudley [1993, Theorem A]: They provide a necessary and sufficient for I_{kr} to be Hilbert-Schmidt in terms of an integral of the metric entropy of the embedding $I_{kr} : \mathcal{H}_k \rightarrow \mathcal{H}_r$. See also Steinwart [2017, Corollary 5.4] for a similar condition.

From Theorem 4.9, it is easy to show that a GP sample path $f \sim \mathcal{GP}(0, k)$ does *not* belong to the corresponding RKHS \mathcal{H}_k with probability 1 if \mathcal{H}_k is infinite dimensional, as summarized in Corollary 4.10 below. This implies that GP samples are “rougher”, or less regular, than RKHS functions (see also Figure 2). Note that this fact has been well known in the literature; see e.g., Wahba [1990, p. 5] and Lukić and Beder [2001, Corollary 7.1].

Corollary 4.10 *Let k be a positive definite kernel on a set \mathcal{X} and \mathcal{H}_k be its RKHS, and consider $f \sim \mathcal{GP}(m, k)$ with $m : \mathcal{X} \rightarrow \mathbb{R}$ satisfying $m \in \mathcal{H}_k$. Then if \mathcal{H}_k is infinite dimensional, then $f \in \mathcal{H}_k$ with probability 0. If \mathcal{H}_k is finite dimensional, then there is a version \tilde{f} of f such that $\tilde{f} \in \mathcal{H}_k$ with probability 1.*

Proof Consider Theorem 4.9 with $r := k$, and let $I_{kk} : \mathcal{H}_k \rightarrow \mathcal{H}_k$ be the inclusion operator, which is the identity map. Let $(\psi_i)_{i \in I} \subset \mathcal{H}_k$ be an orthonormal basis of \mathcal{H}_k , where $|I| = \infty$ if \mathcal{H}_k is infinite dimensional, and $|I| < \infty$ if \mathcal{H}_k is finite dimensional. Then $\|I_{kr}\|_{\text{HS}}^2 = \sum_{i \in I} \|I_{kk}\psi_i\|_{\mathcal{H}_r}^2 = \sum_{i \in I} \|\psi_i\|_{\mathcal{H}_r}^2 = \sum_{i \in I} 1$. Thus, $\|I_{kr}\|_{\text{HS}} = \infty$ if $|I| = \infty$, and $\|I_{kr}\|_{\text{HS}} < \infty$ if $|I| < \infty$. The assertion then follows from Theorem 4.9. \blacksquare

Remark 4.9 Based on the KL expansion (55), Wahba [1990, p. 5] gave an intuitive, but rather heuristic argument to show that a GP sample path does not belong to the corresponding RKHS almost surely; see also Berlinet and Thomas-Agnan [2004, p. 66] and Rasmussen and Williams [2006, Section 6.1]. The argument is as follows. For $f \sim \mathcal{GP}(0, k)$, consider a KL-expansion $f = \sum_{i=1}^{\infty} z_i \lambda_i^{1/2} \phi_i$ with $z_i \sim \mathcal{N}(0, 1)$, where $(\lambda_i^{1/2} \phi_i)_{i=1}^{\infty}$ is an ONB of the RKHS \mathcal{H}_k , which is assumed to be infinite dimensional. Defining $f_m := \sum_{i=1}^m z_i \lambda_i^{1/2} \phi_i$ for $m \in \mathbb{N}$, the KL-expansion may be written as $f = \lim_{m \rightarrow \infty} f_m$. Then,

$$\mathbb{E}[\|f_m\|_{\mathcal{H}_k}^2] = \mathbb{E}\left[\sum_{i=1}^m z_i^2\right] = \sum_{i=1}^m \mathbb{E}[z_i^2] = \sum_{i=1}^m 1 = m.$$

Therefore we have $\lim_{m \rightarrow \infty} \mathbb{E}[\|f_m\|_{\mathcal{H}_k}^2] = \infty$. This *may* imply that $\mathbb{E}[\|f\|_{\mathcal{H}_k}^2] = \infty$, and further that $f \notin \mathcal{H}_k$ with probability 1. Note that, while this argument is intuitive, it is *not* a proof. This is because, as shown in Theorem 4.3, the standard result for the convergence of the KL-expansion $f = \lim_{m \rightarrow \infty} f_m$ is in the mean-square sense (or in $L_2(\nu)$, as mentioned in Remark 4.4). That is, the convergence of the KL-expansion is, of course, *weaker* than the convergence in the RKHS norm, and therefore $\lim_{m \rightarrow \infty} \mathbb{E}[\|f_m\|_{\mathcal{H}_k}^2] = \infty$ does *not* imply $\mathbb{E}[\|f\|_{\mathcal{H}_k}^2] = \infty$. This shows that the importance of carefully considering the convergence type of the KL-expansion, which was investigated and used for establishing GP-sample path properties by Steinwart [2017].

The following example, which follows from Corollary 4.10, recovers the well-known fact that Brownian motion is “non-smooth” while it is continuous.

Example 4.1 *Let f be the standard Brownian motion on $[0, 1]$, which is a Gaussian process with kernel $k(x, y) = \min(x, y)$ for $x, y \in [0, 1]$. The corresponding RKHS \mathcal{H}_k is a Cameron-Martin space [Adler and Taylor, 2007, p. 68] given by*

$$\mathcal{H}_k = \left\{ f \in L_2([0, 1]) : Df \text{ exists and } \int (Df(x))^2 dx < \infty \right\},$$

where Df denotes the weak derivative of f ; this is the first-order Sobolev space on $[0, 1]$. Corollary 4.10 implies that f does not belong to \mathcal{H}_k almost surely. In other words, the Brownian motion does not admit a square-integrable weak derivative.

4.3 Powers of RKHSs as GP Sample Spaces

Driscoll’s theorem (Theorem 4.9) shows a necessary and sufficient condition for a version of $f \sim \mathcal{GP}(m, k)$ to be a member of an RKHS \mathcal{H}_r , but it does not directly provide a way of constructing the RKHS \mathcal{H}_r (nor its reproducing kernel r) based on the given covariance kernel k . This is what is done in Steinwart [2017]: \mathcal{H}_r can be constructed as a *power* of the RKHS \mathcal{H}_k , and r as the corresponding power of the kernel k ; these are concepts introduced by Steinwart and Scovel [2012, Definition 4.1] based on Mercer’s theorem. We review this result, showing that it can be easily derived from Theorem 4.9.

For simplicity, we assume here that a set \mathcal{X} is a compact metric space, a measure ν is a finite Borel measure with \mathcal{X} being its support, and a kernel k is continuous on \mathcal{X} . However, we note that the results of Steinwart [2017] and Steinwart and Scovel [2012] hold under much weaker assumptions (while statements of the results should be modified accordingly). We first introduce the definition of powers of RKHSs and kernels [Steinwart and Scovel, 2012, Definition 4.1].

Definition 4.11 (Powers of RKHSs and kernels) *Let \mathcal{X} be a compact metric space, k be a continuous kernel on \mathcal{X} with \mathcal{H}_k being its RKHS, and ν be a finite Borel measure whose support is \mathcal{X} . Let $0 < \theta \leq 1$ be a constant, and assume that $\sum_{i \in I} \lambda_i^\theta \phi_i^2(x) < \infty$ holds for all $x \in \mathcal{X}$, where $(\lambda_i, \phi_i)_{i \in I}$ is the eigensystem of the integral operator in (47). Then the θ -th power of RKHS \mathcal{H}_k is defined as*

$$\mathcal{H}_k^\theta := \left\{ f = \sum_{i \in I} a_i \lambda_i^{\theta/2} \phi_i : \sum_{i \in I} a_i^2 < \infty \right\}, \quad (57)$$

where the inner-product is given by

$$\langle f, g \rangle_{\mathcal{H}_k^\theta} = \sum_{i \in I} \alpha_i \beta_i \quad \text{for} \quad f := \sum_{i \in I} \alpha_i \lambda_i^{\theta/2} \phi_i \in \mathcal{H}_k, \quad g := \sum_{i \in I} \beta_i \lambda_i^{\theta/2} \phi_i \in \mathcal{H}_k.$$

The θ -th power of kernel k is a function $k^\theta : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ defined by

$$k^\theta(x, y) := \sum_{i \in I} \lambda_i^\theta \phi_i(x) \phi_i(y), \quad x, y \in \mathcal{X}. \quad (58)$$

Remark 4.10 The space \mathcal{H}_k^θ defined as in (57) is in fact an RKHS, with its reproducing kernel being the θ -th power of kernel (58), and \mathcal{H}_k^θ and k^θ are uniquely determined independent of the chosen ONB $(\lambda_i^{1/2}\phi_i)_{i \in I}$ [Steinwart and Scovel, 2012, Proposition 4.2].

Remark 4.11 The power of the RKHS (57) is an intermediate space (or more precisely, an interpolation space) between $L_2(\nu)$ and \mathcal{H}_k , and the constant $0 < \theta \leq 1$ determines how close \mathcal{H}_k^θ is to \mathcal{H}_k [Steinwart and Scovel, 2012, Theorem 4.6]. For instance, if $\theta = 1$ we have $\mathcal{H}_k^\theta = \mathcal{H}_k$, and \mathcal{H}_k^θ approaches $L_2(\nu)$ as $\theta \rightarrow +0$. Indeed, \mathcal{H}_k^θ is nesting with respect to θ :

$$\mathcal{H}_k = \mathcal{H}_k^1 \subset \mathcal{H}_k^\theta \subset \mathcal{H}_k^{\theta'} \subset L_2(\nu), \quad \text{for all } 0 < \theta' < \theta < 1.$$

In other words, \mathcal{H}_k^θ gets larger as θ decreases. If \mathcal{H}_k is an RKHS consisting of smooth functions (such as Sobolev spaces), then \mathcal{H}_k^θ contains less smooth functions than those in \mathcal{H}_k .

The following result, which follows from Theorem 4.9, provides a characterization of GP-sample spaces in terms of powers of RKHSs \mathcal{H}_k^θ . It is a special case of Steinwart [2017, Theorem 5.2], where assumptions required for \mathcal{X} , k and ν are much weaker.

Theorem 4.12 *Let \mathcal{X} be a compact metric space, k be a continuous kernel on \mathcal{X} with \mathcal{H}_k being its RKHS, and ν be a finite Borel measure whose support is \mathcal{X} . Let $0 < \theta < 1$ be a constant, and assume that $\sum_{i \in I} \lambda_i^\theta \phi_i^2(x) < \infty$ holds for all $x \in \mathcal{X}$, where $(\lambda_i, \phi_i)_{i \in I}$ is the eigensystem of the integral operator in (47). Consider $\mathfrak{f} \sim \mathcal{GP}(0, k)$. Then the following statements are equivalent.*

1. $\sum_{i \in I} \lambda_i^{1-\theta} < \infty$.
2. The inclusion operator $I_{kk^\theta} : \mathcal{H}_k \rightarrow \mathcal{H}_k^\theta$ is Hilbert-Schmidt.
3. There exists a version $\tilde{\mathfrak{f}}$ of \mathfrak{f} such that $\tilde{\mathfrak{f}} \in \mathcal{H}_k^\theta$ with probability 1.

Proof The equivalence between 2. and 3. follows from Theorem 4.9 and the fact that \mathcal{H}_k^θ is an RKHS with k^θ being its kernel. The equivalence between 1. and 2. follows from

$$\|I_{kk^\theta}\|_{\text{HS}}^2 = \sum_{i \in I} \|I_{kk^\theta} \lambda_i^{1/2} \phi_i\|_{\mathcal{H}_k^\theta}^2 = \sum_{i \in I} \|\lambda_i^{(1-\theta)/2} \lambda_i^{\theta/2} \phi_i\|_{\mathcal{H}_k^\theta}^2 = \sum_{i \in I} \lambda_i^{1-\theta},$$

where the first equality uses the definition of the Hilbert-Schmidt norm and the fact that $(\lambda_i^{1/2}\phi_i)_{i \in I}$ is an ONB of \mathcal{H}_k , and the third follows from $(\lambda_i^{\theta/2}\phi_i)_{i \in I}$ being an ONB of \mathcal{H}_k^θ . ■

Remark 4.12 Theorem 4.12 shows that the power of the RKHS \mathcal{H}_k^θ contains the support of $\mathfrak{f} \sim \mathcal{GP}(0, k)$, if the eigenvalues $(\lambda_i)_{i \in I}$ satisfy $\sum_{i=1}^\infty \lambda_i^{1-\theta} < \infty$ for $0 < \theta < 1$. Therefore, if one knows the eigensystem $(\lambda_i, \phi_i)_{i \in I}$ of the integral operator (47), one may construct the GP-sample space as \mathcal{H}_k^θ with largest $0 < \theta < 1$ satisfying $\sum_{i=1}^\infty \lambda_i^{1-\theta} < \infty$. Note that the condition $\sum_{i=1}^\infty \lambda_i^{1-\theta} < \infty$ is stronger for larger θ , requiring that the eigenvalues should decay more rapidly (when $|I| = \infty$). Also note that functions in \mathcal{H}_k^θ get smoother as θ increases.

4.4 Examples of Sample Path Properties

We provide concrete examples of GP sample path properties, as corollaries of the above results. We first show sample path properties for GPs with square-exponential kernels in Example 2.1. Intuitively, the result follows from Theorem 4.12 and that the eigenvalues for a square-exponential kernel decay exponentially fast; see Section A.2 for a complete proof.

Corollary 4.13 (Sample path properties for square-exponential kernels) *Let $\mathcal{X} \subset \mathbb{R}^d$ be a compact set with Lipschitz boundary, ν be the Lebesgue measure on \mathcal{X} , k_γ be a square-exponential kernel on \mathcal{X} with bandwidth $\gamma > 0$ with \mathcal{H}_{k_γ} being its RKHS. Then for all $0 < \theta < 1$, the θ -th power $\mathcal{H}_{k_\gamma}^\theta$ of \mathcal{H}_{k_γ} in Definition 4.11 is well-defined. Moreover, for a given $f \sim \mathcal{GP}(0, k_\gamma)$, there exists a version \tilde{f} such that $\tilde{f} \in \mathcal{H}_{k_\gamma}^\theta$ with probability 1 for all $0 < \theta < 1$.*

Remark 4.13 Since $\mathcal{H}_{k_\gamma} \subset \mathcal{H}_{k_\gamma}^\theta$ for $0 < \theta < 1$ and $\mathcal{H}_{k_\gamma}^\theta$ approaches \mathcal{H}_{k_γ} as $\theta \rightarrow 1$, Corollary 4.13 shows that, informally, a GP sample path associated with a square-exponential kernel lies in a space that is infinitesimally larger than \mathcal{H}_{k_γ} . Therefore in practice one should not worry too much about the fact that a GP sample path almost surely falls outside the RKHS \mathcal{H}_{k_γ} , because it nevertheless lies almost surely on the “infinitesimally small shell” surrounding \mathcal{H}_{k_γ} . However, note that the situation is different for Matérn kernels, of which the RKHSs have only a finite degree of smoothness; the above property for square-exponential kernels follows from, intuitively, that functions in the resulting RKHSs are infinitely smooth.

Corollary 4.15 below provides sample path properties for GPs associated with Matérn kernels. To state it, we need to introduce the *interior cone condition* [Wendland, 2005, Definition 3.6], which requires that there is no ‘pinch point’ (i.e. a \prec -shape region) on the boundary of \mathcal{X} .

Definition 4.14 (Interior cone condition) *A set $\mathcal{X} \subset \mathbb{R}^d$ is said to satisfy an interior cone condition if there exist an angle $\theta \in (0, 2\pi)$ and a radius $R > 0$ such that every $x \in \mathcal{X}$ is associated with a unit vector $\xi(x)$ so that the cone $C(x, \xi(x), \psi, R)$ is contained in Ω , where*

$$C(x, \xi(x), \psi, R) := \{x + ay : y \in \mathbb{R}^d, \|y\| = 1, \langle y, \xi(x) \rangle \geq \cos \psi, a \in [0, R]\}.$$

Corollary 4.15 (Sample path properties for Matérn kernels) *Let $\mathcal{X} \subset \mathbb{R}^d$ be a bounded open set such that the boundary is Lipschitz and an interior cone condition is satisfied, and $k_{\alpha, h}$ be the Matérn kernel on \mathcal{X} in Example 2.2 with parameters $\alpha > 0$ and $h > 0$ such that $\alpha + d/2 \in \mathbb{N}$. Then for a given $f \sim \mathcal{GP}(0, k_{\alpha, h})$, there exists a version \tilde{f} such that $\tilde{f} \in \mathcal{H}_{k_{\alpha', h'}}$ with probability 1 for all $\alpha', h' > 0$ satisfying $\alpha > \alpha' + d/2 \in \mathbb{N}$, where $\mathcal{H}_{k_{\alpha', h'}}$ is the RKHS of the Matérn kernel $k_{\alpha', h'}$ with parameters α' and h' .*

Proof Let $s := \alpha + d/2$ and $\beta := \alpha' + d/2$. Denote by $W_2^s(\mathcal{X})$ and $W_2^\beta(\mathcal{X})$ the Sobolev spaces of order s and β respectively, as defined in Example 2.6. Since \mathcal{X} satisfies an interior cone condition and we have $s - \beta > d/2$, Maurin’s theorem [Adams and Fournier, 2003, Theorem 6.61] implies that the embedding $W_2^s(\mathcal{X}) \rightarrow W_2^\beta(\mathcal{X})$ is Hilbert-Schmidt. Since the boundary of \mathcal{X} is Lipschitz, by Wendland [2005, Corollary 10.48] the RKHS $\mathcal{H}_{k_{\alpha, h}}$ of $k_{\alpha, h}$ is norm-equivalent to $W_2^s(\mathcal{X})$, and $\mathcal{H}_{k_{\alpha', h'}}$ is norm-equivalent to $W_2^\beta(\mathcal{X})$. (See also Example 2.6.) Therefore the embedding $\mathcal{H}_{k_{\alpha, h}} \rightarrow \mathcal{H}_{k_{\alpha', h'}}$ is also Hilbert-Schmidt. The assertion then follows from Theorem 4.9. \blacksquare

Remark 4.14 Recall that, as shown in Example 2.6, the RKHS $\mathcal{H}_{k_{\alpha,h}}$ of the Matérn kernel $k_{\alpha,h}$ is norm-equivalent to the Sobolev space $W_2^s(\mathcal{X})$ of order $s := \alpha + d/2$, and $\mathcal{H}_{k_{\alpha',h'}}$ is norm-equivalent to $W_2^\beta(\mathcal{X})$ with $\beta := \alpha' + d/2$. From the assumption $\alpha > \alpha' + d/2$, we have $s > \beta + d/2$. Therefore, Corollary 4.15 shows that, roughly, the smoothness of a GP sample path with $k_{\alpha,h}$, which is β , is $d/2$ -smaller than the smoothness s of the RKHS $\mathcal{H}_{k_{\alpha,h}}$.

Remark 4.15 The condition $\alpha > \alpha' + d/2$ can be shown to be sharp: If $\alpha = \alpha' + d/2$, a sample path $f \sim \mathcal{GP}(0, k_{\alpha,h})$ does *not* belong to $\mathcal{H}_{k_{\alpha',h'}}$ almost surely [Steinwart, 2017, Corollary 5.6, ii]. Also note that the condition $\alpha' + d/2 \in \mathbb{N}$ may be removed; α can be any positive real satisfying $\alpha > \alpha' + d/2$ [Steinwart, 2017, Corollary 5.6, i].

5 Convergence and Posterior Contraction Rates in Regression

In this section, we review asymptotic convergence results for Gaussian process and kernel ridge regression. For both approaches, there have been extensive theoretical investigations, but it seems that the connections between the obtained results for the two approaches are rarely discussed. We therefore discuss the connections between the convergence results for the two approaches in Section 5.1, and show that there is indeed a certain equivalence that highlights the role of regularization and the output noise assumption. The key role in showing this equivalence is played by sample path properties discussed in Section 4. We also review theoretical results from the kernel interpolation literature in Section 5.2. Thanks to the worst case error viewpoint explained in Section 3.4, these results provide upper-bounds for marginal posterior variances in GP-regression. Such bounds are useful in understanding what factors affect the speed of contraction of marginal posterior variances.

5.1 Convergence Rates for Gaussian Process and Kernel Ridge Regression

We here review existing convergence results for Gaussian process regression and kernel ridge regression. Specifically, we compare the posterior contraction rates for GP-regression provided by van der Vaart and van Zanten [2011] and the convergence rates for kernel ridge regression by Steinwart et al. [2009]. The rates obtained in these papers are minimax optimal for regression in Sobolev spaces. Thus, it would be natural to ask how these two results are related. We show that, by focusing on regression in Sobolev spaces, the rates of van der Vaart and van Zanten [2011] can be recovered from those of Steinwart et al. [2009]. This highlights the equivalence between regularization in kernel ridge regression and the additive noise assumption in Gaussian process regression. The arguments are based on sample path properties reviewed in Section 4.

Gaussian process regression. The following model is considered in van der Vaart and van Zanten [2011]. Let $\mathcal{X} = [0, 1]^d$ be the input space, and $f_0 : [0, 1]^d \rightarrow \mathbb{R}$ be the unknown regression function to be estimated. Let $(x, y) \in [0, 1]^d \times \mathbb{R}$ be a joint random variable such that $x \sim P_{\mathcal{X}}$ for a distribution $P_{\mathcal{X}}$ on $[0, 1]^d$ and

$$y = f_0(x) + \varepsilon, \quad \varepsilon \sim \mathcal{N}(0, \sigma^2), \quad (59)$$

where $\sigma^2 > 0$ is the variance of independent noise ε . $P_{\mathcal{X}}$ is assumed to have a density function that is bounded away from zero and infinity. Denote by P the joint distribution of (x, y) , and

assume that one is given an i.i.d. sample $\mathcal{D}_n := (x_i, y_i)_{i=1}^n$ of size n from P as training data. A prior for f_0 is a zero-mean Gaussian process $\mathcal{GP}(0, k_s)$ with covariance kernel k_s , where k_s denotes the Matérn kernel (Example 3) whose RKHS is norm-equivalent to the Sobolev space $W_2^s[0, 1]^d$ of order $s > d/2$. (In the notation of Example 3, this corresponds to $\alpha := s - d/2$; see Example 2.6 for the RKHSs of Matérn kernels.)

GP-regression is performed based on this GP prior, training data \mathcal{D}_n and the likelihood given by the additive Gaussian noise model (59). van der Vaart and van Zanten [2011, Theorem 5] provided the following posterior contraction rate for this setting, assuming the unknown f_0 belongs to a Sobolev space $W_2^\beta[0, 1]^d$ of order $\beta > d/2$. Below $C^\beta([0, 1]^d)$ denotes the Hölder space of order β , and $L_2(P_{\mathcal{X}})$ the Hilbert space of square-integrable functions with respect to $P_{\mathcal{X}}$.

Theorem 5.1 *Let k_s be a kernel on $[0, 1]^d$ whose RKHS is norm-equivalent to the Sobolev space $W_2^s([0, 1]^d)$ of order $s := \alpha + d/2$ with $\alpha > 0$. If $f_0 \in C^\beta([0, 1]^d) \cap W_2^\beta([0, 1]^d)$ and $\min(\alpha, \beta) > d/2$, then we have*

$$\mathbb{E}_{\mathcal{D}_n|f_0} \left[\int \|f - f_0\|_{L_2(P_{\mathcal{X}})}^2 d\Pi_n(f|\mathcal{D}_n) \right] = O(n^{-2\min(\alpha, \beta)/(2\alpha+d)}) \quad (n \rightarrow \infty), \quad (60)$$

where $\mathbb{E}_{\mathcal{X}, Y|f_0}$ denotes the expectation with respect to $\mathcal{D}_n = (x_i, y_i)_{i=1}^n$ with the model (59), and $\Pi_n(f|\mathcal{D}_n)$ the posterior distribution given by GP-regression with kernel k_s .

Remark 5.1 By definition, the posterior mean function (16) is given as $\bar{m}_n := \int f d\Pi_n(f|\mathcal{D}_n)$ (here we made the dependence of m_n on sample size n explicit). It is easy to show that

$$\|\bar{m}_n - f_0\|_{L_2(P_{\mathcal{X}})}^2 \leq \int \|f - f_0\|_{L_2(P_{\mathcal{X}})}^2 d\Pi_n(f|\mathcal{D}_n),$$

and therefore (60) implies the convergence of \bar{m}_n to f_0 ,

$$\mathbb{E}_{\mathcal{D}_n|f_0} \left[\|\bar{m}_n - f_0\|_{L_2(P_{\mathcal{X}})}^2 \right] = O(n^{-2\min(\alpha, \beta)/(2\alpha+d)}) \quad (n \rightarrow \infty). \quad (61)$$

Remark 5.2 The best rate with (61) is attained when $\alpha = \beta$, which results in the rate $n^{-2\beta/(2\beta+d)}$; this is the minimax-optimal rate for regression of a function in $W_2^\beta([0, 1]^d)$ [Stone, 1980, Tsybakov, 2008]. Note that $\alpha = s - d/2$ is essentially the smoothness of $f \sim \mathcal{GP}(0, k_s)$, a sample path of the GP prior (see Corollary 4.15 and Remark 4.14). Therefore the requirement $\alpha = \beta$ means that the smoothness α of sample paths from the GP prior should match the smoothness β of the regression function f_0 .

For later comparison with kernel ridge regression, we point out here that the smoothness s of the corresponding Sobolev RKHS $H^s([0, 1]^d)$ should be specified as $s = \beta + d/2$ to attain the optimal rate. In other words, the smoothness of the RKHS $H^s([0, 1]^d)$ should be *greater* than the Sobolev space $H^\beta([0, 1]^d)$ to which f_0 belongs. This may be counterintuitive, given the equivalence between GP-regression and kernel ridge regression. As we show below, however, the above consequence can be explained from the modeling assumption (59) that noise variance σ^2 remains constant even if n increases.

Kernel ridge regression. For kernel ridge regression, we discuss the convergence results of Steinwart et al. [2009], which do not require the true regression function f_0 be in the RKHS. Let \mathcal{X} be an arbitrary measurable space and $\mathcal{Y} := [-M, M] \subset \mathbb{R}$ be the output space, where $M > 0$ is some constant, and P be a joint distribution on $\mathcal{X} \times \mathcal{Y}$. The unknown regression function $f_0 : \mathcal{X} \rightarrow [-M, M]$ is defined as the conditional expectation $f_0(x) := \mathbb{E}[y|x]$ as usual, where the expectation is with respect to the conditional distribution of y given x with $(x, y) \sim P$. Let $P_{\mathcal{X}}$ be the marginal distribution of P on \mathcal{X} , and assume that it has a density bounded away from zero and infinity; this is the same assumption as for van der Vaart and van Zanten [2011].

Given a training sample $(x_1, y_1), \dots, (x_n, y_n) \stackrel{i.i.d.}{\sim} P$, let \hat{f}_λ be the estimator of f_0 by kernel ridge regression defined as the solution of (23), with $\lambda > 0$ being a regularization constant (here we made the dependence on λ explicit). In Steinwart et al. [2009], the following clipped version \check{f}_λ of \hat{f}_λ is considered for theoretical analysis:

$$\check{f}_\lambda(x) := \begin{cases} M & (\hat{f}_\lambda(x) > M) \\ \hat{f}_\lambda(x) & (-M \leq \hat{f}_\lambda(x) \leq M) \\ -M & (\hat{f}_\lambda(x) < -M) \end{cases}$$

The following result follows from Corollary 6 in Steinwart et al. [2009], the arguments in the paragraphs following Theorem 9 in Steinwart et al. [2009], and the fact $H^\beta([0, 1]^d) \subset B_{2,\infty}^\beta([0, 1]^d)$, where $\beta > d/2$ and $H^\beta([0, 1]^d)$ is the Sobolev space of order $\beta > d/2$ and $B_{2,\infty}^\beta([0, 1]^d)$ is a certain Besov space of the same order β ; see. e.g. Edmunds and Triebel [1996, Eq.7 in p.59].

Theorem 5.2 *Let k_s be a kernel on $\mathcal{X} := [0, 1]^d$ whose RKHS is norm-equivalent to the Sobolev space $W_2^s([0, 1]^d)$ of order $s > d/2$. Assume that $f_0 \in W_2^\beta([0, 1]^d)$ for some $d/2 \leq \beta \leq s$. If $\lambda_n > 0$ is set as*

$$\lambda_n = cn^{-2s/(2\beta+d)}, \quad (62)$$

for a constant $c > 0$ independent of n , then we have

$$\|\check{f}_{\lambda_n} - f_0\|_{L_2(P_{\mathcal{X}})}^2 = O_p(n^{-2\beta/(2\beta+d)}) \quad (n \rightarrow \infty). \quad (63)$$

Remark 5.3 As mentioned earlier, the rate (63) is minimax optimal for regression in the Sobolev space $W_2^\beta([0, 1]^d)$ of order β [Stone, 1980, Tsybakov, 2008].

Remark 5.4 Theorem 5.2 does not require that f_0 be in the RKHS of kernel k_s : the smoothness β of f_0 can be smaller than the smoothness s of the RKHS, in which case f_0 does not belong to $H^s([0, 1]^d)$. Intuitively, this is possible because a function outside the RKHS but “not very far away” from the RKHS can be approximated well by functions in the RKHS. This intuition is in fact formally characterized and exploited in Steinwart et al. [2009] by using approximation theory based on interpolation spaces. Note also that the degree of approximation is controlled by the regularization schedule (62): λ_n should decrease more quickly, as the smoothness β of f_0 becomes smaller.

The following corollary is a special case of Theorem 5.2, which is essentially equivalent to Theorem 5.1 for GP-regression.

Corollary 5.3 *Assume that $f_0 \in W_2^\beta([0, 1]^d)$ for $\beta > 0$. Let k_s be a kernel on $\mathcal{X} := [0, 1]^d$ whose RKHS is norm-equivalent to the Sobolev space $W_2^s([0, 1]^d)$ of order $s := \beta + d/2$, and define $\lambda_n := cn^{-1}$ with $c > 0$ being any constant. Then (63) holds for \check{f}_{λ_n} .*

Remark 5.5 Recall that the variance σ^2 of output noise in GP-regression is related to the regularization constant λ in kernel ridge regression as $\sigma^2 = n\lambda_n$. Therefore, the modeling assumption that σ^2 does not vary with n in GP-regression is equivalent to the regularization schedule $\lambda_n = cn^{-1}$ in kernel ridge regression. With this regularization schedule, the optimal rate (63) is attained by kernel k_s with $s = \beta + d/2$. This optimal order s of the kernel is the same as for the optimal order in Theorem 5.1 for GP-regression (i.e., $s = \alpha + d/2$ with $\alpha = \beta$). Thus, Corollary 5.3 is essentially equivalent to Theorem 5.1, revealing a theoretical equivalence between GP and kernel ridge regression.

5.2 Upper-bounds and Contraction Rates for Posterior Variance

We focus here on the noise-free case, and review the results that provide contraction rates for posterior variance $\bar{k}(x, x)$ in GP-interpolation (22). It seems that these results have not been well known in the machine learning literature. However, we believe that they are important in particular in understanding the mechanism of active learning methods based on GPs and Bayesian optimization, as these methods make use of the posterior variance function in exploring new points to evaluate. In fact, these results have essentially been used in Bull [2011] for theoretical analysis of Bayesian optimization; see also e.g., Briol et al. [2018], Tuo and Wu [2016], Stuart and Teckentrup [2018] for similar applications of these results.

The results we review are from the literature on kernel interpolation [Wendland, 2005, Schaback and Wendland, 2006, Scheuere et al., 2013]. In the kernel interpolation literature, the posterior standard deviation function $(\bar{k}(x, x))^{1/2}$ is called the *power function*, and has been studied because it provides an upper-bound for the error of kernel interpolation, as can be seen from Corollary 3.11. The key role is played by the quantity called the *fill distance*, which quantifies the denseness of points $X = \{x_1, \dots, x_n\} \subset \mathbb{R}^d$ in the region of interest $\mathcal{X} \subset \mathbb{R}^d$. For a constant $\rho > 0$, the fill distance at $x \in \mathcal{X}$ is defined by

$$h_{\rho, X}(x) := \sup_{y \in \mathcal{X}: \|x-y\| \leq \rho} \min_{x_i \in X} \|y - x_i\|. \quad (64)$$

In other words, by thinking of the ball $B(x, \rho)$ around x of radius ρ , the fill distance $h_{\rho, X}(x)$ is the radius of the maximum ball in $B(x, \rho)$ in which no points are contained. Thus, $h_{\rho, X}(x)$ being smaller implies that more points are located around x .

The following result, which is from Wu and Schaback [1993, Theorem 5.14], provides an upper-bound for the posterior variance in terms of the fill distance (64). In particular it applies to cases where the kernel induces an RKHS that is norm-equivalent to Sobolev spaces (e.g., Matérn kernels).

Theorem 5.4 *Let k be a kernel on \mathbb{R}^d whose RKHS is norm-equivalent to the Sobolev space of order s . Then for any $\rho > 0$, there exist constants $h_0 > 0$ and $C > 0$ satisfying the following: For any $x \in \mathbb{R}^d$ and any set of points $X = \{x_1, \dots, x_n\} \subset \mathbb{R}^d$ satisfying $h_{\rho, X}(x) \leq h_0$, we have*

$$\bar{k}(x, x) \leq Ch_{\rho, X}^{2s-d}(x). \quad (65)$$

Remark 5.6 For a kernel with infinite smoothness (such as square-exponential kernels), the exponent in the upper-bound (65) can be arbitrarily large. That is, for $\rho > 0$ and any $\alpha > 0$, there exist constants $h_\alpha > 0$ and $C_\alpha > 0$ satisfying the following: For any $x \in \mathbb{R}^d$ and any set of points $X = \{x_1, \dots, x_n\} \subset \mathbb{R}^d$ satisfying $h_{\rho, X}(x) \leq h_\alpha$, we have

$$\bar{k}(x, x) \leq C_\alpha h_{\rho, X}^\alpha(x).$$

Note that constants C_α and h_α depend on α , so the upper-bound may not monotonically decrease as α increases, for fixed points X .

Remark 5.7 An upper-bound of exponential order can be derived for a kernel with infinite smoothness, but this technically requires that the fill distance be defined globally on the region of interest $\mathcal{X} \subset \mathbb{R}^d$: For $X = \{x_1, \dots, x_n\} \subset \mathcal{X}$, the (global) fill distance $h_X > 0$ is defined as

$$h_X := \sup_{x \in \mathcal{X}} \min_{x_i \in X} \|x - x_i\|.$$

For instance, if \mathcal{X} is a cube in \mathbb{R}^d and k is a square-exponential kernel on \mathcal{X} , Wendland [2005, Theorem 11.22] shows that there exists a constant $c > 0$ that does not depend on h_X such that

$$\bar{k}(x, x) \leq \exp(c \log(h_{X, \Omega}) / \sqrt{h_X}).$$

whenever h_X is sufficiently small.

Theorem 5.4 shows that the amount of posterior variance $\bar{k}(x, x)$ is determined by i) the local fill distance $h_{\rho, X}$, ii) the dimensionality d of the space, and iii) the smoothness s of the kernel. This implies that $\bar{k}(x, x)$ contracts to 0 as the denseness of the points x_1, \dots, x_n around x increases, and that the rate of contraction is determined by d and s . This is formally characterized by the following corollary, which directly follows from Theorem 5.4.

Corollary 5.5 *Let k be a kernel on \mathbb{R}^d whose RKHS is norm-equivalent to the Sobolev space of order $s > d/2$, and fix $\rho > 0$. For $x \in \mathbb{R}^d$, assume that $X = (x_1, \dots, x_n) \subset \mathbb{R}^d$ satisfy $h_{\rho, X}(x) = O(n^{-b})$ as $n \rightarrow \infty$ for some $b > 0$. Then we have*

$$\bar{k}(x, x) = O(n^{-b(2s-d)}) \quad (n \rightarrow \infty). \tag{66}$$

Remark 5.8 If x_1, \dots, x_n are given as equally-spaced grid points in the ball of radius ρ around x , then it can be easily shown that $h_{\rho, X}(x) = O(n^{-1/d})$ as $n \rightarrow \infty$. Thus in this case, the rate in (66) becomes $\bar{k}(x, x) = O(n^{-(2s/d-1)})$, which reveals the existence of the curse of dimensionality: the required number of points increases exponentially in the dimension d to achieve a certain level of posterior contraction.

6 Integral Transforms

This section reviews some examples of the connections between RKHS and GP approaches that involve integrals of the kernel. Such computations arise both in the context of estimating a latent (probability) measure and when estimating the integral of a latent function against a known measure. In Section 6.1, we first introduce kernel mean embeddings of distributions

and resulting metrics on probability measures (Maximum Mean Discrepancy), and then provide their probabilistic interpretations based on Gaussian processes. We next describe their connections to kernel and Bayesian quadrature, which are approaches to numerical integration based on positive definite kernels in Section 6.2. Coming back to the statistical setting, a kernel mean shrinkage estimator and its Bayesian interpretation are discussed in Section 6.3. Finally, we review a nonparametric dependency measure called Hilbert-Schmidt Independence Criterion, which is defined via kernel mean embeddings, and present its probabilistic interpretation based on Gaussian processes.

In this section, we will use the following notation to denote integrals:

$$Pf := \int f(x)dP(x), \quad P_n f := \sum_{i=1}^n w_i f(x_i),$$

where P is a measure on a measurable space \mathcal{X} , $P_n := \sum_{i=1}^n w_i \delta_{x_i}$ is an empirical measure on \mathcal{X} with $(w_i, x_i)_{i=1}^n \subset \mathbb{R} \times \mathcal{X}$ and δ_{x_i} being a Dirac distribution at x_i , and $f : \mathcal{X} \rightarrow \mathbb{R}$ is a function.

6.1 Maximum Mean Discrepancy: Worst Case and Average Case Errors

Let $(\mathcal{X}, \mathfrak{B})$ be a measurable space, k be a bounded measurable kernel on \mathcal{X} with \mathcal{H}_k being its RKHS, and \mathcal{P} be the set of all probability measures on \mathcal{X} . For any $P \in \mathcal{P}$, its *kernel mean* (or mean embedding) is defined as the integral of the canonical feature map $k(\cdot, x)$ with respect to P :

$$\mu_P := \int k(\cdot, x)dP(x) \in \mathcal{H}_k, \quad (67)$$

which is well-defined as a Bochner integral, as long as $\int \sqrt{k(x, x)}dP(x) < \infty$ [Sriperumbudur et al., 2010, Theorem 1]. The kernel mean (67) is an element in RKHS \mathcal{H}_k that represents P .

Remark 6.1 The notion of embeddings of probability measures can be readily extended to embeddings of finite signed measures [Sriperumbudur et al., 2011]. This is important because in the case where P is an empirical approximation or estimator of the form $P_n := \sum_{i=1}^n w_i \delta_{x_i}$, the weights w_1, \dots, w_n can be negative. For instance, this is the case of kernel or Bayesian quadrature discussed in Section 6.2.

Characteristic kernels and MMD. If the mapping $P \in \mathcal{P} \mapsto \mu_P \in \mathcal{H}_k$ is injective, that is, if $\mu_P = \mu_Q$ implies $P = Q$ for any probability measures P and Q on \mathcal{X} , then the kernel k is called *characteristic* [Fukumizu et al., 2004, 2008, Sriperumbudur et al., 2010]. For instance, characteristic kernels on $\mathcal{X} = \mathbb{R}^d$ include the Gaussian and Matérn kernels [Sriperumbudur et al., 2010]. If kernel k is characteristic, each kernel mean μ_P is uniquely associated with the embedded measure P , and therefore one can define a distance between probability measures P and Q as the distance between the kernel means μ_P and μ_Q in the RKHS:

$$\text{MMD}(P, Q; \mathcal{H}_k) := \|\mu_P - \mu_Q\|_{\mathcal{H}_k} = \sup_{\substack{f \in \mathcal{H}_k: \\ \|f\|_{\mathcal{H}_k} \leq 1}} (Pf - Qf) = \sup_{\substack{f \in \mathcal{H}_k: \\ \|f\|_{\mathcal{H}_k} \leq 1}} |Pf - Qf|, \quad (68)$$

where the second equality follows from \mathcal{H}_k being a vector space and the Cauchy-Schwartz inequality and the third follows from \mathcal{H}_k being a vector space; see the proof of Gretton et al.

[2012, Lemma 4]. Because of this expression, this distance between kernel means is called *maximum mean discrepancy* (MMD) [Gretton et al., 2012], as it is the maximum difference between integrals (means) Pf and Qf , when f is taken from the unit ball in RKHS \mathcal{H}_k . If k is characteristic, MMD becomes a proper metric on probability measures, and thus its estimator can be used as a test statistic in hypothesis testing for the two sample problem or goodness of fit [Gretton et al., 2012, Chwialkowski et al., 2016, Liu et al., 2016].

Remark 6.2 Let k be a bounded, continuous shift-invariant kernel on $\mathcal{X} = \mathbb{R}^d$ such that $k(x, y) = \phi(x - y)$ for $x, y \in \mathbb{R}^d$, where $\phi : \mathbb{R}^d \rightarrow \mathbb{R}$ is a positive definite function. For such a kernel, Sriperumbudur et al. [2010, Corollary 4] provided a spectral characterization of MMD in terms of the Fourier transform $\mathcal{F}[\Phi]$ of Φ : For any Borel probability measures P and Q on \mathbb{R}^d , it holds that

$$\text{MMD}^2(P, Q, \mathcal{H}_k) = \int |\psi_P(\omega) - \psi_Q(\omega)|^2 \mathcal{F}[\Phi](\omega) d\omega$$

where ψ_P and ψ_Q are the characteristic functions of P and Q , respectively. That is, MMD between P and Q is the weighted L_2 distance between the characteristic functions ψ_P and ψ_Q , where the weight function is the Fourier transform $\mathcal{F}[\phi]$. From this expression and the fact that characteristic functions and distributions are one-to-one, Sriperumbudur et al. [2010] provided a necessary and sufficient condition for a shift-invariant kernel to be characteristic: A bounded continuous shift invariant kernel k is characteristic if and only if the support of the Fourier transform $\mathcal{F}[\phi]$ is \mathbb{R}^d [Sriperumbudur et al., 2010, Theorem 9].

Remark 6.3 The use of RKHS in defining (68) is practically convenient because, thanks to the reproducing property, the squared MMD can be written as

$$\|\mu_P - \mu_Q\|_{\mathcal{H}_k}^2 = \mathbb{E}_{x, x'}[k(x, x')] - 2\mathbb{E}_{x, y}[k(x, y)] + \mathbb{E}_{y, y'}[k(y, y')],$$

where $x, x' \sim P$ and $y, y' \sim Q$ are all independent [Gretton et al., 2012, Lemma 6]. Therefore, by replacing the expectations in the right side by empirical ones, one can straightforwardly estimate the MMD from samples.

Worst case error. In the literature on numerical integration or quasi Monte Carlo, the right side of (68) is known as the *worst case error* [Hickernell, 1998, Dick et al., 2013]. To be more precise, in the problem of numerical integration or sampling, P is a *known* probability measure and $Q := P_n := \sum_{i=1}^n w_i \delta_{x_i}$ is its approximation, where $(w_i, x_i)_{i=1}^n \subset \mathbb{R} \times \mathcal{X}$ are generated by a user so that P_n becomes an accurate approximation of P . As such, one is interested in the quality of approximation of P_n to P . For this purpose, (68) is used as a quantitative measure of approximation, and is interpreted as the worst case error of numerical integration $|Pf - P_n f|$ when f is taken from the unit ball in RKHS \mathcal{H}_k . We will discuss this problem of numerical integration in detail in Section 6.2.

Probabilistic interpretation as the average case error. Proposition 6.1 below provides a probabilistic interpretation of MMD in terms of the GP of the kernel k . More specifically, MMD between P and Q can be understood as the *expected squared difference* between integrals Pf and Qf , where the expectation is with respect to a draw f from $\mathcal{GP}(0, k)$. In the terminology

of numerical integration, this shows the equivalence between the RKHS worst case error and the Gaussian process *average case error*. While this equivalence has been known in the literature [Ritter, 2000, Corollary 7 in p.40], we provide a proof, as it is instructive.

Proposition 6.1 *Let k be a bounded kernel on a measurable space \mathcal{X} , and P and Q be finite measures on \mathcal{X} . Then we have*

$$\text{MMD}^2(P, Q; \mathcal{H}_k) = \left(\sup_{\|f\|_{\mathcal{H}_k} \leq 1} (Pf - Qf) \right)^2 = \mathbb{E}_{f \sim \mathcal{GP}(0, k)} \left[(Pf - Qf)^2 \right]. \quad (69)$$

Proof Let x and x' be independent random variables following P , y and y' be those following Q , and f be an independent draw from $\mathcal{GP}(0, k)$. By using the reproducing property and the expression (5) of the kernel $k(x, y) = \mathbb{E}_f[f(x)f(y)]$

$$\begin{aligned} \|\mu_P - \mu_Q\|_{\mathcal{H}_k}^2 &= \mathbb{E}_{x, x'}[k(x, x')] - 2\mathbb{E}_{x, y}[k(x, y)] + \mathbb{E}_{y, y'}[k(y, y')] \\ &= \mathbb{E}_{x, x'}[\mathbb{E}_f f(x)f(x')] - 2\mathbb{E}_{x, y}[\mathbb{E}_f f(x)f(y)] + \mathbb{E}_{y, y'}[\mathbb{E}_f f(y)f(y')] \\ &\stackrel{*}{=} \mathbb{E}_f \left[\mathbb{E}_{x, x'}[f(x)f(x')] - 2\mathbb{E}_{x, y}[f(x)f(y)] + \mathbb{E}_{y, y'}[f(y)f(y')] \right] \\ &= \mathbb{E}_f \left[(\mathbb{E}_x[f(x)])^2 - 2\mathbb{E}_x[f(x)]\mathbb{E}_y[f(y)] + (\mathbb{E}_y[f(y)])^2 \right] \\ &= \mathbb{E}_f \left[(\mathbb{E}_x[f(x)] - \mathbb{E}_y[f(y)])^2 \right], \end{aligned}$$

where $\stackrel{*}{=}$ follows from Fubini's theorem, which is applicable because we have

$$\mathbb{E}_{x, x'} \mathbb{E}_f |f(x)f(x')| \leq \mathbb{E}_{x, x'} \sqrt{\mathbb{E}_f f^2(x)} \sqrt{\mathbb{E}_f f^2(x')} = \mathbb{E}_x \sqrt{k(x, x)} \mathbb{E}_{x'} \sqrt{k(x', x')} < \infty,$$

where the last inequality follows from k being bounded. ■

Remark 6.4 Since $f \sim \mathcal{GP}(0, k)$ is a mean-zero Gaussian process, the expectation of the real-valued random variable $Pf - Qf$ is 0. Therefore the right side of (69) can be seen as the variance of this random variable $Pf - Qf$:

$$\mathbb{E}_{f \sim \mathcal{GP}(0, k)} \left[(Pf - Qf)^2 \right] = \text{var}[Pf - Qf].$$

Since $Pf - Qf$ is a linear transform of Gaussian process f , it is a real-valued Gaussian random variable. This implies that, when dealing with MMD between P and Q , one implicitly deals with the distribution of $Pf - Qf$, which is $\mathcal{N}(0, \sigma^2)$ where $\sigma^2 := \text{MMD}(P, Q; \mathcal{H}_k)$.

The following corollary immediately follows from Proposition 6.1 and the definition of a kernel being characteristic. It provides a probabilistic interpretation of a characteristic kernel in terms of the corresponding Gaussian process.

Corollary 6.2 *Let k be a bounded characteristic kernel on a measurable space \mathcal{X} . Then for any probability measures P and Q on \mathcal{X} , we have $P = Q$ if and only if $Pf = Qf$ holds almost surely for a Gaussian process $f \sim \mathcal{GP}(0, k)$.*

6.2 Sampling and Numerical Integration

Here we consider the problem of numerical integration or (deterministic) sampling. Let $f : \mathcal{X} \rightarrow \mathbb{R}$ be an integrand and P be a *known* probability measure on \mathcal{X} , and assume that the integral $\int f(x)dP(x)$ cannot be computed analytically. The task is to numerically compute the integral as a weighted sum of function values

$$\sum_{i=1}^n w_i f(x_i) \approx \int f(x)dP(x).$$

Therefore the problem is how to select the weighted points $(w_i, x_i)_{i=1}^n \subset \mathbb{R} \times \mathcal{X}$ so that this approximation becomes as accurate as possible. If one has prior knowledge about certain properties of f such as its smoothness, then one can use this information for the construction of a quadrature rule. This can be done by making use of positive definite kernels.

Kernel quadrature. For simplicity, assume that design points $x_1, \dots, x_n \in \mathcal{X}$ are already given and fixed. In kernel quadrature, weights w_1, \dots, w_n are obtained by the minimization of MMD between $Q := P_n := \sum_i w_i \delta_{x_i}$ and P :

$$\min_{w_1, \dots, w_n \in \mathbb{R}} \text{MMD}(P_n, P; \mathcal{H}_k) = \min_{w_1, \dots, w_n \in \mathbb{R}} \sup_{\|f\|_{\mathcal{H}_k} \leq 1} |P_n f - P f|,$$

which is, as shown in the right side, equivalent to the minimization of the worst case error in the unit ball of RKHS \mathcal{H}_k . Assuming that kernel matrix k_{XX} is invertible, this optimization problem can be solved in closed form, and the resulting weights are given by

$$(w_1, \dots, w_n)^T = k_{XX}^{-1} \mu_X \in \mathbb{R}^n, \tag{70}$$

where $\mu_X := (\mu_P(x_i))_{i=1}^n \in \mathbb{R}^n$. Using these weights, integral Pf is approximated as

$$P_n f = \sum_{i=1}^n w_i f(x_i). \tag{71}$$

We assumed here that design points x_1, \dots, x_n are given at the beginning, but there are also approaches that obtain design points by aiming at the minimization of the worst case error; examples include quasi Monte Carlo methods [Dick et al., 2013] and kernel herding [Chen et al., 2010].

Remark 6.5 To calculate the weights in (70), one needs to be able to evaluate function values of the kernel mean $\mu_P = \int k(\cdot, x)dP(x)$, thus requiring a certain compatibility between k and P . For instance, this is possible when k is square-exponential and P is Gaussian on $\mathcal{X} \subset \mathbb{R}^d$. For other examples, see Briol et al. [2018, Table 1]. Note that one is also able to perform kernel quadrature using kernel Stein discrepancy; in this case the weights (70) can be calculated if one knows the gradient of log density of P [Oates et al., 2017, Liu and Lee, 2017]. This remark also applies to Bayesian quadrature explained below.

Bayesian quadrature. Bayesian quadrature [Diaconis, 1988, O’Hagan, 1991, Briol et al., 2018, Karvonen et al., 2018] is a probabilistic approach to numerical integration based on Gaussian processes. As before, assume for simplicity that design points x_1, \dots, x_n are fixed. In this approach, a prior distribution is put on the integrand f as a Gaussian process $f \sim \mathcal{GP}(0, k)$. Assume that functions values $f(x_1), \dots, f(x_n)$ at the design points are provided. In Bayesian quadrature, these input-output pairs $(x_i, f(x_i))_{i=1}^n$ are regarded as “observed data.” Then the posterior distribution of the integral is given by

$$Pf \mid (x_i, f(x_i))_{i=1}^n \sim \mathcal{N}(\mu_n, \sigma_n^2), \quad (72)$$

where $\mu_n \in \mathbb{R}$ and $\sigma_n^2 > 0$ are respectively the posterior mean and variance given by

$$\mu_n := \mu_X^\top k_{XX}^{-1} f_X = P_n f, \quad (73)$$

$$\sigma_n^2 := \int \int k(x, x') dP(x) dP(x') - \mu_X^\top k_{XX}^{-1} \mu_X, \quad (74)$$

where $\mu_X := (\mu_P(x_i))_{i=1}^n \in \mathbb{R}^n$ with μ_P being the kernel mean and $P_n := \sum_{i=1}^n w_i \delta_{x_i}$.

Remark 6.6 As discussed in Section 3.1, since we deal with noise-free observations $f(x_1), \dots, f(x_n)$, there is no likelihood model in the above derivation (or, the likelihood function is degenerate). Therefore (72) is not a “posterior distribution” in the usual sense of Bayesian inference. However, (72) is still well-defined as a conditional distribution of Pf given $(x_i, f(x_i))_{i=1}^n$ [Cockayne et al., 2017, Section 2.5]. For discussion regarding what it means by “Bayesian” in the noise-free setting or in numerical analysis, we refer to Cockayne et al. [2017].

Remark 6.7 Note that the posterior variance (73) does not depend on the given integrand f , and determined only by the kernel k , the design points x_1, \dots, x_n and the measure P .

First note that the posterior mean (73) is identical to the integral estimate (71) of kernel quadrature. The following result shows that the posterior variance (74) of Bayesian quadrature is also equal to the squared MMD between P_n and P , where the weights w_1, \dots, w_n are given in (70). This identity has been known in the literature [Huszár and Duvenaud, 2012, Briol et al., 2018], but we provide a proof, as it is simple and instructive.

Proposition 6.3 *Let σ_n^2 be the posterior variance (74) of Bayesian quadrature, and $P_n := \sum_{i=1}^n w_i \delta_{x_i}$ be the empirical measure with the weights w_1, \dots, w_n given in (70). Then*

$$\sigma_n^2 = \|\mu_{P_n} - \mu_P\|_{\mathcal{H}_k}^2 = \text{MMD}^2(P_n, P; \mathcal{H}_k). \quad (75)$$

Proof By the reproducing property and the definition of the weights $w = (w_1, \dots, w_n)^\top \in \mathbb{R}^n$, we have

$$\begin{aligned} \|\mu_{P_n} - \mu_P\|_{\mathcal{H}_k}^2 &= \|\mu_{P_n}\|_{\mathcal{H}_k}^2 - 2 \langle \mu_{P_n}, \mu_P \rangle_{\mathcal{H}_k} + \|\mu_P\|_{\mathcal{H}_k}^2 \\ &= w^\top k_{XX} w - 2w^\top \mu_X + \int \int k(x, x') dP(x) dP(x') \\ &= -\mu_X^\top k_{XX}^{-1} \mu_X + \int \int k(x, x') dP(x) dP(x'), \end{aligned}$$

and the result follows. ■

Remark 6.8 Proposition 6.3 shows the equivalence between the average case error w.r.t. GPs and the (squared) worst case error in the RKHS, in the setting of numerical integration. This is because (75) can be written as

$$\mathbb{E}_f [(Pf - \mu_n)^2 \mid (x_i, f(x_i))_{i=1}^n] = \left(\sup_{\|f\|_{\mathcal{H}_k} \leq 1} |P_n f - Pf| \right)^2.$$

This equivalence has been used by Briol et al. [2015, 2018] to establish posterior contraction rates of Bayesian quadrature, by transferring results on the worst case error [Bach et al., 2012, Dick et al., 2013] to the probabilistic or Bayesian setting.

Noisy observations and robustness to misspecification. If one’s knowledge about integrand f of interest is limited, it could happen that it does not belong to RKHS \mathcal{H}_k , that is, misspecification of the hypothesis class may occur. Kanagawa et al. [2016, 2017] showed that kernel quadrature can be made robust to such misspecification, by introducing a quadratic regularizer for the weights, i.e.,

$$\min_{w_1, \dots, w_n \in \mathbb{R}^n} \text{MMD}^2(P_n, P; \mathcal{H}_k) + \lambda \sum_{i=1}^n w_i^2, \quad (76)$$

where $\lambda > 0$ is a regularization constant. The resulting weights are then given by

$$(w_1, \dots, w_n)^\top = (k_{XX} + n\lambda I_n)^{-1} \mu_X \in \mathbb{R}^n. \quad (77)$$

Kernel quadrature with quadratic weight regularization has also been studied by Bach [2017], who pointed out such regularization is equivalent to assuming the existence of additive noises in the function values. That is, assume that, instead of observing the exact function values $f(x_1), \dots, f(x_n)$, one is given noisy observations $y_i = f(x_i) + \varepsilon_i$, where $\varepsilon_i \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \sigma^2)$ with $\sigma^2 := n\lambda$. Bayesian quadrature under this assumption yields the posterior distribution of the integral as $Pf \mid (x_i, y_i)_{i=1}^n \sim \mathcal{N}(\mu_n, \sigma_n^2)$, where

$$\begin{aligned} \mu_n &:= \mu_X^\top (k_{XX} + \sigma^2 I_n)^{-1} Y = \sum_{i=1}^n w_i y_i, \\ \sigma_n^2 &:= \int \int k(x, x') dP(x) dP(x') - \mu_X^\top (k_{XX} + \sigma^2 I_n)^{-1} \mu_X, \end{aligned}$$

where w_1, \dots, w_n are given by (77). These are regularized versions of (73) and (74), with f_X replaced by $Y = (y_1, \dots, y_n)^\top \in \mathbb{R}^n$.

Discussion on the difference between the kernel and Bayesian approaches. The optimization viewpoint of kernel quadrature allows one to directly incorporate a constraint on quadrature weights w_1, \dots, w_n . For instance, Liu and Lee [2017] proposed to optimize the weights under a constraint that the weights be non-negative. Such a constraint is not straightforward to be realized only with a probabilistic perspective. On the other hand, with Bayesian quadrature one can express prior knowledge about the integrand that is not easy to be incorporated with the kernel approach. For instance, Gunter et al. [2014] proposed to model an integrand that is non-negative as a squared GP (i.e., as a chi-squared process). Such modeling can be realized because one expresses the prior knowledge as a generative model, but this is not straightforward to achieve through the optimization of weights.

6.3 Kernel Mean Shrinkage Estimator and Its Bayesian Interpretation

Given an i.i.d. sample $x_1, \dots, x_n \sim P$, an empirical estimator of the kernel mean (67) is defined as

$$\hat{\mu}_P := \frac{1}{n} \sum_{i=1}^n k(\cdot, x_i), \quad (78)$$

which satisfies $\mathbb{E}[\|\hat{\mu}_P - \mu_P\|_{\mathcal{H}_k}] = O(n^{-1/2})$ as $n \rightarrow \infty$, if k is bounded [Smola et al., 2007, Tolstikhin et al., 2017]. One way to compute MMD empirically is to substitute this estimate (and that for μ_Q) in (68): this results in a V-statistic estimate of MMD [Gretton et al., 2012, Eq. 5]. Tolstikhin et al. [2017] showed that the rate $n^{-1/2}$ is minimax-optimal and thus cannot be improved, meaning that (78) is an *asymptotically* optimal estimator. However, when the sample size n is *fixed*, it is known that there exists a “better” estimator that improves upon (78), which was shown by Muandet et al. [2016]. More precisely, consider an estimator $\hat{\mu}_{P,\alpha}$ defined by

$$\hat{\mu}_{P,\alpha} := f^* + (1 - \alpha)\hat{\mu}_P, \quad (79)$$

where $f^* \in \mathcal{H}_k$ is fixed and arbitrary and α is a constant. Muandet et al. [2016, Theorem 1] proved that, if (and only if) the constant satisfies $0 < \alpha < 2\mathbb{E}[\|\hat{\mu}_P - \mu_P\|_{\mathcal{H}_k}^2]/(\mathbb{E}[\|\hat{\mu}_P - \mu_P\|_{\mathcal{H}_k}^2] + \|f^* - \mu_P\|_{\mathcal{H}_k}^2)$, then $\hat{\mu}_{P,\alpha}$ produces a smaller mean-squared error than $\hat{\mu}_P$:

$$\mathbb{E}[\|\hat{\mu}_{P,\alpha} - \mu_P\|_{\mathcal{H}_k}^2] < \mathbb{E}[\|\hat{\mu}_P - \mu_P\|_{\mathcal{H}_k}^2].$$

A motivation for Muandet et al. [2016] was the so called *Stein phenomenon*, which states that the standard empirical estimator for the mean of a d -dimensional Gaussian distribution with $d \geq 3$ is *inadmissible*, meaning that there exists an estimator that yields smaller mean squared error for a fixed sample size [Stein, 1956]. Motivated by this old result, Muandet et al. [2016] proposed a number of shrinkage estimators including the one in (79), some of which have been theoretically proven to be “better” than the standard empirical estimator in (78) in terms of the mean squared RKHS error.

We review here a certain shrinkage estimator proposed by Muandet et al. [2016, Section 4] called the *spectral kernel mean estimator* (SKME), and its Bayesian interpretation given by Flaxman et al. [2016]; this provides another instance of the connection between the kernel and Bayesian approaches. Different from (79), however, the SKME has *not* been shown to be theoretically better than the empirical estimator (78), but has only been shown to yield better empirical performance. Given an i.i.d. sample $X = (x_1, \dots, x_n) \sim P$, the SKME is defined as

$$\check{\mu}_{P,\lambda} := \sum_{i=1}^n w_i k(\cdot, x_i) \quad (80)$$

where the weights $w_1, \dots, w_n \in \mathbb{R}$ are given by

$$(w_1, \dots, w_n)^T := (k_{XX} + n\lambda I_n)^{-1} \hat{\mu}_X \in \mathbb{R}^n, \quad (81)$$

with $\hat{\mu}_X := (\hat{\mu}_P(x_i))_{i=1}^n \in \mathbb{R}^n$, $k_{XX} \in \mathbb{R}^{n \times n}$ being the kernel matrix, and $\lambda > 0$ being a regularization constant. The estimator in (80) was originally derived from a certain RKHS-valued ridge regression problem. Alternatively, the estimator can be derived by solving the

following minimization problem:

$$\min_{w_1, \dots, w_n \in \mathbb{R}} \left\| \sum_{i=1}^n w_i k(\cdot, X_i) - \hat{\mu}_P \right\|_{\mathcal{H}}^2 + \lambda \|w\|^2, \quad (82)$$

where $\hat{\mu}_P$ is the empirical estimator in (78). The weights in (81) are given as the solution of this optimization problem. This interpretation shows that as the regularization constant λ increases, the estimator (80) shrinks towards zero in the RKHS, thus reducing the variance of the estimator while introducing bias. Therefore λ controls the bias-variance trade off; this is beneficial in practice when the sample size is relatively small, in which case the variance of the empirical estimator (78) may be large.

Remark 6.9 The weights (81) of the SKME are essentially the same as those for kernel quadrature with quadratic regularization (77). The only difference is that, while the information of the true kernel mean μ_P is used in (77), the empirical mean $\hat{\mu}_P$ is used in (81), since in the statistical setting μ_P is an unknown quantity to be estimated. Note also the essential equivalence of the two optimization problems (76) and (82), based on which the weights (77) and (81) are respectively derived.

Bayesian interpretation of the shrinkage estimator. Flaxman et al. [2016] showed that there exists a Bayesian interpretation of the shrinkage estimator in (80). We formulate their approach by using the powered kernel (58) in Section 4.3. The prior for the kernel mean μ_P is defined as a Gaussian process:

$$\mu_P \sim \mathcal{GP}(0, k^\theta), \quad (83)$$

where k^θ is the powered kernel (58) with the power $\theta \geq 1$ appropriately chosen so that μ_P can be a sample path of the GP. On the other hand, by regarding the evaluations $\hat{\mu}_P$ at sample points x_1, \dots, x_n as ‘‘observations’’, Flaxman et al. [2016] proposed to define a likelihood function as an additive Gaussian-noise model

$$\hat{\mu}_P(x_i) = \mu_P(x_i) + \varepsilon_i, \quad \varepsilon_i \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \sigma^2), \quad i = 1, \dots, n. \quad (84)$$

Note that these assumptions are essentially those of GP-regression. (We will discuss the validity of this assumption shortly.) Therefore a similar argument as in Section 3.1 implies that the posterior distribution of μ_P is also a GP and given by

$$\mu_P \mid (x_i, \hat{\mu}_P(x_i))_{i=1}^n \sim \mathcal{GP}(\bar{\mu}_P, \bar{k}^\theta),$$

where $\bar{\mu}_P$ and \bar{k}^θ are the posterior mean function and the posterior covariance function, respectively, and are given by

$$\bar{\mu}_P(x) := k_{xX}^\theta (k_{XX}^\theta + \sigma^2 I_n)^{-1} \hat{\mu}_X, \quad x \in \mathcal{X}, \quad (85)$$

$$\bar{k}^\theta(x, x') := k^\theta(x, x') - k_{xX}^\theta (k_{XX}^\theta + \sigma^2 I_n)^{-1} k_{Xx'}^\theta, \quad x, x' \in \mathcal{X}, \quad (86)$$

where $k_{Xx}^\theta = k_{xX}^{\theta \top} := (k^\theta(x, x_i))_{i=1}^n \in \mathbb{R}^n$, $k_{XX}^\theta := (k^\theta(x_i, x_j)) \in \mathbb{R}^{n \times n}$ and $\hat{\mu}_X := (\hat{\mu}_P(x_i))_{i=1}^n \in \mathbb{R}^n$.

Remark 6.10 If $\theta = 1$ and $\lambda = \sigma^2/n$, the posterior mean function (85) is equal to the shrinkage estimator (80). Therefore in this case, the modeling assumptions (83) and (84) provide a probabilistic interpretation for the shrinkage estimator (80). In other words, the shrinkage estimator (80) implicitly performs Bayesian inference under the probabilistic model (83) and (84). This probabilistic viewpoint turns out to be practically useful. For instance, Flaxman et al. [2016] made use of the probabilistic formulation for selecting the kernel hyperparameter (and the noise variance σ^2) in an unsupervised fashion, by using the empirical Bayes method; Law et al. [2018] used the posterior covariance function (86) to enable uncertainty quantification in application to distribution regression.

Discussion on the modeling assumption. We make some remarks on the modeling assumptions (83) and (84). First, as mentioned above, the GP prior (83) should be defined so that the kernel mean μ_P can be a sample path of the GP. If $\theta = 1$, this is not the case: As reviewed earlier, GP sample paths do not belong to the RKHS of the covariance kernel with probability one. This fact motivated Flaxman et al. [2016] to use a certain kernel that is smoother than the kernel defining the kernel mean, in order to guarantee that the kernel mean can be a GP sample path. We instead defined the GP prior (83) using the powered kernel k^θ : If $\theta > 1$, the kernel k^θ is “smoother” than the original kernel k , and there is “sufficiently large” θ to guarantee that a GP sample path lies in the RKHS. For instance, if k is a square-exponential kernel it can be shown from Corollary 4.13 that, the choice $\theta = 1 + \varepsilon$ with $\varepsilon > 0$ being arbitrarily small guarantees that sample paths of $\mathcal{GP}(0, k^\theta)$ belong to the RKHS of k with probability one. In other words, k^θ can be chosen so that the resulting power of RKHS \mathcal{H}_k^θ is “infinitesimally smaller” than the original RKHS \mathcal{H}_k . This may explain why the use of $\theta = 1$ with a square-exponential kernel resulted in good empirical performance in Law et al. [2018].

We note that Flaxman et al. [2016] introduced the likelihood model (84) for computational feasibility, i.e., to obtain the posterior distribution of μ_P as a GP. Therefore, the assumption (84) may not be conceptually well-motivated. For instance, if the kernel k is bounded, so is $\hat{\mu}_P = \frac{1}{n} \sum_{i=1}^n k(\cdot, x_i)$; this fact is not captured by the assumption that noise ε_i is additive Gaussian, which is unbounded. Moreover, noise ε_i is neither independent nor of constant variance in general, since the differences $\hat{\mu}_P(x_i) - \mu_P(x_i)$ at different points x_i are dependent.

Open questions. If $\theta > 1$, which is required to ensure that μ_P is a sample path of $\mathcal{GP}(0, k^\theta)$, then the resulting posterior mean function (85) does not coincide with the shrinkage estimator (80). This raises the following question: Can the use of the smoother kernel k^θ lead to “better” performance for estimation of $\mu_P = \int k(\cdot, x)dP(x)$ in terms of the mean-square error with a fixed sample size, when compared to the standard empirical estimator (78)? Muandet et al. [2016] was not able to show such superiority of the shrinkage estimator; this may be because they used the kernel k^θ with $\theta = 1$ in (80), which is not supported from the Bayesian interpretation.

Relation to Bayesian quadrature. For any RKHS function $f \in \mathcal{H}$, the integral $\int f(x)dP(x)$ can be estimated as a weighted sum $\sum_{i=1}^n w_i f(x_i)$, where the weighted points $(w_i, x_i)_{i=1}^n$ are

those expressing $\check{\mu}_{P,\lambda}$ as in (80). This follows from the inequality

$$\left| \sum_{i=1}^n w_i f(x_i) - \int f(x) dP(x) \right| = \left| \langle \check{\mu}_{P,\lambda} - \mu_P, f \rangle_{\mathcal{H}_k} \right| \leq \| \check{\mu}_{P,\lambda} - \mu_P \|_{\mathcal{H}_k} \| f \|_{\mathcal{H}_k}$$

and that $\check{\mu}_P$ should be close to μ_P . It is easy to show that the weighted sum can be written as

$$\sum_{i=1}^n w_i f(x_i) = \frac{1}{n} \sum_{i=1}^n \bar{m}(x_i),$$

where $\bar{m} : \mathcal{X} \rightarrow \mathbb{R}$ is the posterior mean function (16) in GP regression. In other words, the weighted sum is equal to the empirical mean of the the fitted function \bar{m} . As shown in Section 6.2, this is essentially Bayesian quadrature using the empirical measure $P_n = \frac{1}{n} \sum_{i=1}^n \delta_{x_i}$.

6.4 Gaussian Process Interpretation of Hilbert Schmidt Independence Criterion

Here we deal with a popular kernel-based dependency measure known as the *Hilbert-Schmidt Independence Criterion* (HSIC) [Gretton et al., 2005], which can be defined as the MMD between the joint distribution of two random variables and the product of their marginals. HSIC is a nonparametric dependency measure, and as such does not require a specific parametric assumption about the form of dependencies between random variable variables. Since it also can be calculated in a simple way using kernels, it has found a wide range of applications including independence testing [Gretton et al., 2008, Zhang et al., 2018], variable selection [Song et al., 2012, Yamada et al., 2014], post selection inference [Yamada et al., 2018], and causal discovery [Pfister et al., 2017], to name a few. We provide here a probabilistic interpretation of HSIC in terms on GPs; to the best of our knowledge, this probabilistic interpretation of general HSIC measures is novel and it recovers Brownian distance covariance of Székely and Rizzo [2009], known to be a special case of HSIC [Sejdinovic et al., 2013].

Let X and Y be random variables taking values in measurable spaces \mathcal{X} and \mathcal{Y} respectively. Denote by $P_{\mathcal{X} \times \mathcal{Y}}$ the joint distribution of X and Y , and let $P_{\mathcal{X}}$ and $P_{\mathcal{Y}}$ be its marginal distributions on \mathcal{X} and \mathcal{Y} , respectively. Let k and ℓ be positive definite kernels on \mathcal{X} and \mathcal{Y} respectively, and let $\mathcal{H}_{\mathcal{X}}$ and $\mathcal{H}_{\mathcal{Y}}$ be their respective RKHSs. For the product kernel $k \otimes \ell : (\mathcal{X} \times \mathcal{Y}) \times (\mathcal{X} \times \mathcal{Y}) \rightarrow \mathbb{R}$ defined by $(k \otimes \ell)((x, y), (x', y')) := k(x, x')\ell(y, y')$, the corresponding RKHS is denoted by $\mathcal{H}_{\mathcal{X}} \otimes \mathcal{H}_{\mathcal{Y}}$, which is the tensor product of $\mathcal{H}_{\mathcal{X}}$ and $\mathcal{H}_{\mathcal{Y}}$.

While HSIC was originally defined by Gretton et al. [2005] as the Hilbert-Schmidt norm of a certain cross-covariance operator [Fukumizu et al., 2004], we follow here an equivalent definition given by Smola et al. [2007]: HSIC is defined as the (squared) MMD between $P_{\mathcal{X} \times \mathcal{Y}}$ and $P_{\mathcal{X}}P_{\mathcal{Y}}$ in the RKHS $\mathcal{H}_{\mathcal{X}} \otimes \mathcal{H}_{\mathcal{Y}}$:

$$\text{HSIC}(X, Y) := \text{MMD}^2(P_{\mathcal{X} \times \mathcal{Y}}, P_{\mathcal{X}}P_{\mathcal{Y}}) = \| \mu_{P_{\mathcal{X} \times \mathcal{Y}}} - \mu_{P_{\mathcal{X}}} \otimes \mu_{P_{\mathcal{Y}}} \|_{\mathcal{H}_{\mathcal{X}} \otimes \mathcal{H}_{\mathcal{Y}}}^2, \quad (87)$$

where $\mu_{P_{\mathcal{X} \times \mathcal{Y}}}$ and $\mu_{P_{\mathcal{X}}} \otimes \mu_{P_{\mathcal{Y}}}$ are the kernel means of $P_{\mathcal{X} \times \mathcal{Y}}$ and $P_{\mathcal{X}}P_{\mathcal{Y}}$, respectively. If the kernel $k \otimes \ell$ is characteristic, then the HSIC is zero if and only if X and Y are independent; see Szabó and Sriperumbudur [2018] for thorough analysis of conditions for $k \otimes \ell$ being characteristic. Thus in this case, HSIC is qualified as a nonparametric measure of dependence.

Thanks to the reproducing property, HSIC can be expressed in terms of expectations of the kernels [Gretton et al., 2005, Lemma 1]:

$$\begin{aligned} \text{HSIC}(X, Y, k, \ell) &= \mathbb{E}_{X, Y, X', Y'} [k(X, X')\ell(Y, Y')] \\ &\quad - 2\mathbb{E}_{X, Y} [\mathbb{E}_{X'}[k(X, X')]\mathbb{E}_{Y'}\ell(Y, Y')] + \mathbb{E}_{X, X'} [k(X, X')] \mathbb{E}_{Y, Y'} [\ell(Y, Y')], \end{aligned} \quad (88)$$

where X' and Y' are respectively independent copies of X and Y . Given an i.i.d. sample $((X_i, Y_i))_{i=1}^n$ from the joint distribution $P_{\mathcal{X} \times \mathcal{Y}}$, an empirical estimator of HSIC is straightforwardly given by replacing the expectations in (88) by the corresponding empirical averages; for details see Gretton et al. [2005].

Gaussian Process Interpretation. Consider independent draws from the zero-mean GPs of the covariance kernels k and ℓ :

$$\mathbf{f} \sim \mathcal{GP}(0, k), \quad \mathbf{g} \sim \mathcal{GP}(0, \ell).$$

The following result provides a probabilistic interpretation of HSIC in terms of these GPs.

Proposition 6.4 *Let k and ℓ be positive definite kernels on measurable spaces \mathcal{X} and \mathcal{Y} respectively, and let $\mathbf{f} \sim \mathcal{GP}(0, k)$ and $\mathbf{g} \sim \mathcal{GP}(0, \ell)$ be independent Gaussian processes. For random variables X and Y taking values respectively in \mathcal{X} and \mathcal{Y} , we have*

$$\text{HSIC}(X, Y) = \mathbb{E}_{\mathbf{f}, \mathbf{g}} \text{cov}^2(\mathbf{f}(X), \mathbf{g}(Y)), \quad (89)$$

where $\text{HSIC}(X, Y)$ is defined by (87), and

$$\text{cov}(\mathbf{f}(X), \mathbf{g}(Y)) := \mathbb{E}_{X, Y} [(\mathbf{f}(X) - \mathbb{E}[\mathbf{f}(X)|\mathbf{f}]) (\mathbf{g}(Y) - \mathbb{E}[\mathbf{g}(Y)|\mathbf{g}]) | \mathbf{f}, \mathbf{g}].$$

Note that $\text{cov}(\mathbf{f}(X), \mathbf{g}(Y))$ is the covariance between the *real-valued* random variables $\mathbf{f}(X)$ and $\mathbf{g}(Y)$, with \mathbf{f} and \mathbf{g} being fixed. Proposition 6.4 thus shows that HSIC is the expectation of the square of this covariance with respect to the draws $\mathbf{f} \sim \mathcal{GP}(0, k)$ and $\mathbf{g} \sim \mathcal{GP}(0, \ell)$. In other words, the computation of $\text{HSIC}(X, Y)$ amounts to simultaneously considering various nonlinear transformations $\mathbf{f}(X)$ and $\mathbf{g}(Y)$ of random variables X and Y (as defined by the GPs), and then computing the average of the (squared) covariance between these transformed variables. We believe that this interpretation provides a simple way to understand HSIC as a measure of dependence, in particular for people who are familiar with GPs but not with RKHSs.

Connection to Brownian Covariance. We mention that the expression in the right side of (89) is related to the *Brownian (distance) covariance* introduced by Székely and Rizzo [2009]. To describe this, let $X_{\mathbf{f}}$ and $Y_{\mathbf{g}}$ be *real-valued* random variables such that

$$X_{\mathbf{f}} := \mathbf{f}(X) - \mathbb{E}[\mathbf{f}(X)|\mathbf{f}] = \mathbf{f}(X) - \int \mathbf{f}(x) dP_{\mathcal{X}}(x), \quad (90)$$

$$Y_{\mathbf{g}} := \mathbf{g}(Y) - \mathbb{E}[\mathbf{g}(Y)|\mathbf{g}] = \mathbf{g}(Y) - \int \mathbf{g}(y) dP_{\mathcal{Y}}(y). \quad (91)$$

Then the covariance between $\mathbf{f}(X)$ and $\mathbf{g}(Y)$ (with \mathbf{f} and \mathbf{g} being fixed) can be written as

$$\begin{aligned} \text{cov}(\mathbf{f}(X), \mathbf{g}(Y)) &= \mathbb{E}_{X, Y} [(\mathbf{f}(X) - \mathbb{E}[\mathbf{f}(X)|\mathbf{f}]) (\mathbf{g}(Y) - \mathbb{E}[\mathbf{g}(Y)|\mathbf{g}]) | \mathbf{f}, \mathbf{g}] \\ &= \mathbb{E}_{X, Y} [X_{\mathbf{f}} Y_{\mathbf{g}} | \mathbf{f}, \mathbf{g}], \end{aligned}$$

and therefore it follows that

$$\mathbb{E}_{\mathbf{f}, \mathbf{g}} \text{cov}^2(\mathbf{f}(X), \mathbf{g}(Y)) = \mathbb{E}_{\mathbf{f}, \mathbf{g}, (X, Y), (X', Y')} [X_{\mathbf{f}} X'_{\mathbf{f}} Y_{\mathbf{g}} Y'_{\mathbf{g}}], \quad (92)$$

where (X', Y') is an independent copy of the joint random variable (X, Y) , and $X'_{\mathbf{f}}$ and $Y'_{\mathbf{g}}$ are defined similarly to (90) and (91). The right side in (92) coincides with the definition of a dependence measure given by Székely and Rizzo [2009, Definition 5], where they consider as \mathbf{f} and \mathbf{g} arbitrary stochastic processes on Euclidean spaces. Specifically, the right side in (92) is the definition of the Brownian covariance [Székely and Rizzo, 2009, Definition 4], if $\mathcal{X} = \mathbb{R}^p$ and $\mathcal{Y} = \mathbb{R}^q$ for $p, q \in \mathbb{N}$ and if \mathbf{f} and \mathbf{g} are respectively the Brownian motions with the covariance kernels k and ℓ given by

$$k(x, x') := \|x\| + \|x'\| - 2\|x - x'\|, \quad x, x' \in \mathbb{R}^p, \quad (93)$$

$$\ell(y, y') := \|y\| + \|y'\| - 2\|y - y'\|, \quad y, y' \in \mathbb{R}^q. \quad (94)$$

In this case, the Brownian covariance is further identical to the *distance covariance* [Székely and Rizzo, 2009, Theorem 8], a nonparametric measure of dependence for random variables taking values in Euclidean spaces [Székely and Rizzo, 2009, Definition 1]. Therefore our result implies that HSIC is identical to the distance covariance, when the kernels are given by (93) and (94); we thus have recovered the result of Sejdinovic et al. [2013, Theorem 24] based on the probabilistic interpretation of HSIC.

7 Conclusions

In machine learning, statistics and numerical analysis, both the notion of a kernel and that of a Gaussian process play central roles in theoretical analysis. In fact, they are so central in machine learning that they may be seen as placeholders for statistical learning theory and Bayesian analysis, the two mathematical frameworks that have historically provided the theoretical foundation of the field. Kernel methods are founded on notions like regularization and optimization, while Gaussian processes are generative models operating in terms of marginal and conditional distributions. The present text provided a review of the intersection of these two areas, covering both fundamental equivalences and differences. It is important to clarify and understand these relationships to facilitate the transfer of knowledge and methods from one side to the other. At a time when machine learning is arguably expecting the emergence of a third, still only vaguely discernible new theoretical foundation in particular for deep models, this paper is also an opportunity to note that “frequentist” and “Bayesian” statistics are not always as different from each other as they may appear at first sight. We hope that this contribution is an important step towards developing a common language between the two fields, which will lead to further advances in each field, which otherwise would have been much more difficult to achieve.

Acknowledgements

We would like to thank Mark van der Wilk for fruitful discussions. The original idea for this manuscript arose during Workshop 16481 of the Leibniz-Centre for Computer Science at Schloß Dagstuhl. The authors would like to express the Centre for their hospitality and support. MK and PH acknowledge support by the European Research Council (StG Project

PANAMA). BKS is supported by NSF-DMS-1713011. DS is supported in part by The Alan Turing Institute (EP/N510129/1) and by the ERC (FP7/617071). All authors except the first are arranged in an alphabetical order.

A Proofs

A.1 Proof of Lemma 3.9

Proof By the reproducing property, the right side of (38) can be written as

$$\sup_{\|f\|_{\mathcal{H}_k} \leq 1} \sum_{i=1}^m c_i f(x_i) = \sup_{\|f\|_{\mathcal{H}_k} \leq 1} \left\langle \sum_{i=1}^m c_i k(\cdot, x_i), f \right\rangle_{\mathcal{H}_k}. \quad (95)$$

By the Cauchy-Schwartz inequality, the right side of this equality is upper-bounded as

$$\sup_{\|f\|_{\mathcal{H}_k} \leq 1} \left\langle \sum_{i=1}^m c_i k(\cdot, x_i), f \right\rangle_{\mathcal{H}_k} \leq \sup_{\|f\|_{\mathcal{H}_k} \leq 1} \left\| \sum_{i=1}^m c_i k(\cdot, x_i) \right\|_{\mathcal{H}_k} \|f\|_{\mathcal{H}_k} = \left\| \sum_{i=1}^m c_i k(\cdot, x_i) \right\|_{\mathcal{H}_k}.$$

On the other hand, defining $g := \sum_{i=1}^m c_i k(\cdot, x_i) / \left\| \sum_{i=1}^m c_i k(\cdot, x_i) \right\|_{\mathcal{H}_k}$, we have $\|g\|_{\mathcal{H}_k} = 1$, and thus the right side of (95) can be lower-bounded as

$$\sup_{\|f\|_{\mathcal{H}_k} \leq 1} \left\langle \sum_{i=1}^m c_i k(\cdot, x_i), f \right\rangle_{\mathcal{H}_k} \geq \left\langle \sum_{i=1}^m c_i k(\cdot, x_i), g \right\rangle_{\mathcal{H}_k} = \left\| \sum_{i=1}^m c_i k(\cdot, x_i) \right\|_{\mathcal{H}_k}.$$

The assertion follows from (95) and these lower and upper bounds. ■

A.2 Proof of Corollary 4.13

For Banach spaces A and B , we denote by $A \hookrightarrow B$ that $A \subset B$ and that the inclusion is continuous. We first need to the notion of *interpolation spaces*; for details, see e.g., Adams and Fournier [2003, Section 7.6], Steinwart and Christmann [2008, Section 5.6], Cucker and Zhou [2007, Section 4.5] and references therein.

Definition A.1 Let E and F be Banach spaces such that $E \hookrightarrow F$. Let $K : E \times \mathbb{R}_+ \rightarrow \mathbb{R}_+$ be the K -functional defined by

$$K(x, t) := K(x, t, E, F) := \inf_{y \in F} (\|x - y\|_E + t\|y\|_F), \quad x \in E, \quad t > 0.$$

Then for $0 < \theta \leq 1$, the interpolation space $[E, F]_{\theta, 2}$ is a Banach space defined by

$$[E, F]_{\theta, 2} := \left\{ x \in E : \|x\|_{[E, F]_{\theta, 2}} < \infty \right\},$$

where the norm is defined by

$$\|x\|_{[E, F]_{\theta, 2}}^2 := \int_0^\infty \left(t^{-\theta} K(x, t) \right)^2 \frac{dt}{t}.$$

We will need the following lemma.

Lemma A.2 *Let E , F and G be Banach spaces such that $G \hookrightarrow F \hookrightarrow E$. Then we have $[E, G]_{\theta,2} \hookrightarrow [E, F]_{\theta,2}$ for all $0 < \theta \leq 1$.*

Proof First note that, since $G \hookrightarrow F$, there exists a constant $c > 0$ such that $\|z\|_F \leq c\|z\|_G$ holds for all $z \in G$. Therefore, for all $x \in E$ and $t > 0$, we have

$$\begin{aligned} K(x, t, E, F) &= \inf_{y \in F} (\|x - y\|_E + t\|y\|_F) \\ &\leq \inf_{z \in G} (\|x - z\|_E + t\|z\|_F) \\ &\leq \inf_{z \in G} (\|x - z\|_E + ct\|z\|_G) = K(x, ct, E, G). \end{aligned}$$

Thus, we have

$$\begin{aligned} \|x\|_{[E,F]_{\theta,2}}^2 &= \int_0^\infty \left(t^{-\theta} K(x, t, E, F) \right)^2 \frac{dt}{t} \\ &\leq \int_0^\infty \left(t^{-\theta} K(x, ct, E, G) \right)^2 \frac{dt}{t} \\ &= c^{-2\theta} \int_0^\infty \left(s^{-\theta} K(x, s, E, G) \right)^2 \frac{ds}{s} \quad (s := ct) \\ &= c^{-2\theta} \|x\|_{[E,G]_{\theta,2}}^2, \quad x \in [E, G]_{\theta,2}. \end{aligned}$$

which implies the assertion. ■

We are now ready to prove Corollary 4.13.

Proof Fix $\theta \in (0, 1)$. First we show that $\sum_{i=1}^\infty \lambda_i^\theta \phi_i^2(x) < \infty$ holds for all $x \in \mathcal{X}$, which implies that the power of RKHS $\mathcal{H}_{k_\gamma}^\theta$ is well defined. To this end, by Steinwart [2017, Theorem 2.5], it is sufficient to show that

$$[L_2(\nu), \mathcal{H}_{k_\gamma}]_{\theta,2} \hookrightarrow L_\infty(\nu), \tag{96}$$

where $L_2(\nu)$ and $L_\infty(\nu)$ are to be understood as quotient spaces with respect to ν , and \mathcal{H}_{k_γ} as the embedding in $L_2(\nu)$; see Steinwart [2017, Section 2] for precise definition.

Let $m \in \mathbb{N}$ be such that $\theta m > d/2$, and $W_2^m(\mathcal{X})$ be the Sobolev space of order m on \mathcal{X} . By Steinwart and Christmann [2008, Theorem 4.48], we have $\mathcal{H}_{k_\gamma} \hookrightarrow W_2^m(\mathcal{X})$. Therefore Lemma A.2 implies that $[L_2(\nu), \mathcal{H}_{k_\gamma}]_{\theta,2} \hookrightarrow [L_2(\nu), W_2^m(\mathcal{X})]_{\theta,2}$. Note that, since ν is the Lebesgue measure, $[L_2(\nu), W_2^m(\mathcal{X})]_{\theta,2}$ is the Besov space $B_{22}^{\theta m}(\mathcal{X})$ of order θm [Adams and Fournier, 2003, Section 7.32]. Since \mathcal{X} is a bounded Lipschitz domain and $\theta m > d/2$, we have $B_{22}^{\theta m}(\mathcal{X}) \hookrightarrow L_\infty(\nu)$ [Triebel, 2006, Proposition 4.6]; see also Adams and Fournier [2003, Theorem 7.34]. Combining these embeddings implies (96).

We next show that $\sum_{i=1}^\infty \lambda_i^{1-\theta} < \infty$, which implies the assertion by Theorem 4.12. To this end, for $i \in \mathbb{N}$, define the i -th (dyadic) entropy number of the embedding $\text{id} : \mathcal{H}_{k_\gamma} \rightarrow L_2(\nu)$ by

$$\varepsilon_i(\text{id} : \mathcal{H}_{k_\gamma} \rightarrow L_2(\nu)) := \inf \left\{ \varepsilon > 0 : \exists (h_j)_{j=1}^{2^{i-1}} \subset L_2(\nu) \text{ s.t. } B_{\mathcal{H}_{k_\gamma}} \subset \bigcup_{j=1}^{2^{i-1}} (h_j + \varepsilon B_{L_2(\nu)}) \right\},$$

where $B_{\mathcal{H}_{k_\gamma}}$ and $B_{L_2(\nu)}$ denote the centered unit balls in \mathcal{H}_{k_γ} and $L_2(\nu)$, respectively. Note that since \mathcal{X} is bounded, there exists a ball of radius $r \geq \gamma$ that contains \mathcal{X} . From Meister and Steinwart [2016, Theorem 12], for all $p \in (0, 1)$, we have

$$\varepsilon_i(\text{id} : \mathcal{H}_{k_\gamma}(\mathcal{X}) \rightarrow L_2(\nu)) \leq c_{p,d,r,\gamma} i^{-1/2p}, \quad i \geq 1.$$

$c_{p,d,r,\gamma} > 0$ is a constant depending only on p, d, r and γ . Using this inequality, the i -th largest eigenvalue λ_i is upper-bounded as

$$\lambda_i \leq 4\varepsilon_i^2(\text{id} : \mathcal{H}_{k_\gamma} \rightarrow L_2(\nu)) \leq 4c_{p,d,r,\gamma}^2 i^{-1/p}, \quad i \geq 1,$$

where the first inequality follows from Steinwart [2017, Lemma 2.6 Eq. 23]; this lemma is applicable since we have $\int_{\mathcal{X}} k_\gamma(x, x) d\nu(x) < \infty$ and thus the embedding $\text{id} : \mathcal{H}_{k_\gamma} \rightarrow L_2(\nu)$ is compact [Steinwart and Scovel, 2012, Lemma 2.3]. Therefore we have

$$\sum_{i=1}^{\infty} \lambda_i^{1-\theta} < (4c_{p,d,r,\gamma}^2)^{1-\theta} \sum_{i=1}^{\infty} i^{-(1-\theta)/p}.$$

The right side is bounded, if we take $p \in (0, 1)$ such that $1-\theta > p$. This implies $\sum_{i=1}^{\infty} \lambda_i^{1-\theta} < \infty$. \blacksquare

A.3 Proof of Proposition 6.4

Proof Since we have the identity (92), it is sufficient to prove that the right side of (92) is equal to the HSIC (87). First note that

$$\begin{aligned} X_f X'_f &= \left(f(X) - \int f(x) dP_{\mathcal{X}}(x) \right) \left(f(X') - \int f(x) dP_{\mathcal{X}}(x) \right) \\ &= f(X)f(X') - \int f(x)f(X') dP_{\mathcal{X}}(x) - \int f(X)f(x) dP_{\mathcal{X}}(x) + \int \int f(x)f(x') dP_{\mathcal{X}}(x) dP_{\mathcal{X}}(x'). \end{aligned}$$

By using the identity $k(x, x') = \mathbb{E}_f[f(x)f(x')]$ for $x, x' \in \mathcal{X}$, we then obtain

$$\begin{aligned} &\mathbb{E}_f[X_f X'_f | X, X', Y, Y'] \\ &= k(X, X') - \int k(x, X') dP_{\mathcal{X}}(x) - \int k(X, x) dP_{\mathcal{X}}(x) + \int \int k(x, x') dP_{\mathcal{X}}(x) dP_{\mathcal{X}}(x') \\ &= k(X, X') - \mu_{P_{\mathcal{X}}}(X') - \mu_{P_{\mathcal{X}}}(X) + \|\mu_{P_{\mathcal{X}}}\|_{\mathcal{H}_{\mathcal{X}}}^2 \\ &= \langle k(\cdot, X) - \mu_{P_{\mathcal{X}}}, k(\cdot, X') - \mu_{P_{\mathcal{X}}} \rangle_{\mathcal{H}_{\mathcal{X}}}. \end{aligned}$$

Similarly, one can show that

$$\mathbb{E}_f[Y_g Y'_g | X, X', Y, Y'] = \langle \ell(\cdot, Y) - \mu_{P_Y}, \ell(\cdot, Y') - \mu_{P_Y} \rangle_{\mathcal{H}_Y}.$$

Define $\bar{k}(\cdot, X) := k(\cdot, X) - \mu_{P_X}$ and $\bar{l}(\cdot, Y) := l(\cdot, Y) - \mu_{P_Y}$. Therefore, the right side of (92) can be written as

$$\begin{aligned}
& \mathbb{E}[X_f X_f' Y_g Y_g'] \\
&= \mathbb{E}_{X, X', Y, Y'} [\mathbb{E}_f [X_f X_f' | X, X', Y, Y'] \mathbb{E}_g [Y_g Y_g' | X, X', Y, Y']] \\
&= \mathbb{E}_{X, X', Y, Y'} \left[\langle k(\cdot, X) - \mu_{P_X}, k(\cdot, X') - \mu_{P_X} \rangle_{\mathcal{H}_X} \langle \ell(\cdot, Y) - \mu_{P_Y}, \ell(\cdot, Y') - \mu_{P_Y} \rangle_{\mathcal{H}_Y} \right] \\
&= \mathbb{E}_{X, X', Y, Y'} \left[\langle \bar{k}(\cdot, X) \otimes \bar{l}(\cdot, Y), \bar{k}(\cdot, X') \otimes \bar{l}(\cdot, Y') \rangle_{\mathcal{H}_X \otimes \mathcal{H}_Y} \right] \\
&= \langle \mathbb{E}_{X, Y} [\bar{k}(\cdot, X) \otimes \bar{l}(\cdot, Y)], \mathbb{E}_{X', Y'} [\bar{k}(\cdot, X') \otimes \bar{l}(\cdot, Y')] \rangle_{\mathcal{H}_X \otimes \mathcal{H}_Y} \\
&= \left\| \mathbb{E}_{X, Y} [\bar{k}(\cdot, X) \otimes \bar{l}(\cdot, Y)] \right\|_{\mathcal{H}_X \otimes \mathcal{H}_Y}^2 \\
&= \left\| \mathbb{E}_{X, Y} [k(\cdot, X) \otimes \ell(\cdot, Y) - k(\cdot, X) \otimes \mu_{P_Y} - \mu_{P_X} \otimes \ell(\cdot, Y) + \mu_{P_X} \otimes \mu_{P_Y}] \right\|_{\mathcal{H}_X \otimes \mathcal{H}_Y}^2 \\
&= \left\| \mu_{P_{X \times Y}} - \mu_{P_X} \otimes \mu_{P_Y} \right\|_{\mathcal{H}_X \otimes \mathcal{H}_Y}^2 = \text{HSIC}(X, Y).
\end{aligned}$$

■

References

- R. A. Adams and J. J. F. Fournier. *Sobolev Spaces*. Academic Press, New York, 2nd edition, 2003.
- R. J. Adler. *An Introduction to Continuity, Extrema, and Related Topics for General Gaussian Processes*, volume 12. Institute of Mathematical Statistics, 1990.
- R. J. Adler and J. E. Taylor. *Random Fields and Geometry*. Springer, 2007.
- N. Aronszajn. Theory of reproducing kernels. *Transactions of the American Mathematical Society*, 68(3), pages 337–404, 1950.
- F. Bach. On the equivalence between kernel quadrature rules and random feature expansions. *Journal of Machine Learning Research*, 18(19):1–38, 2017.
- F. Bach, S. Lacoste-Julien, and G. Obozinski. On the equivalence between herding and conditional gradient algorithms. In *Proceedings of the 29th International Conference on Machine Learning (ICML2012)*, pages 1359–1366, 2012.
- A. Berlinet and C. Thomas-Agnan. *Reproducing Kernel Hilbert Spaces in Probability and Statistics*. Kluwer, 2004.
- P. Brémaud. *Fourier Analysis and Stochastic Processes*. Springer, 2014.
- F.-X. Briol, C. J. Oates, M. Girolami, and M. A. Osborne. Frank-Wolfe Bayesian quadrature: Probabilistic integration with theoretical guarantees. In *Advances in Neural Information Processing Systems 28*, pages 1162–1170, 2015.

- F.-X. Briol, C. J. Oates, M. Girolami, M. A. Osborne, and D. Sejdinovic. Probabilistic integration: A role in statistical computation? *Statistical Science (to appear)*, *arXiv:1512.00933 [stat.ML]*, 2018.
- A. D. Bull. Convergence rates of efficient global optimization algorithms. *Journal of Machine Learning Research*, 12:2879–2904, 2011.
- A. Caponnetto and E. D. Vito. Optimal rates for regularized least-squares algorithm. *Foundations of Computational Mathematics*, 7(4):331–368, 2007.
- Y. Chen, M. Welling, and A. Smola. Supersamples from kernel-herding. In P. Grünwald and P. Spirtes, editors, *Proceedings of the 26th Conference on Uncertainty in Artificial Intelligence (UAI 2010)*, pages 109–116. AUAI Press, 2010.
- K. Chwialkowski, H. Strathmann, and A. Gretton. A kernel test of goodness of fit. In M. F. Balcan and K. Q. Weinberger, editors, *Proceedings of The 33rd International Conference on Machine Learning*, volume 48 of *Proceedings of Machine Learning Research*, pages 2606–2615. PMLR, 2016.
- J. Cockayne, C. Oates, T. Sullivan, and M. Girolami. Bayesian probabilistic numerical methods. *ArXiv e-prints*, arXiv:1702.03673v2 [stat.ME], Feb. 2017.
- F. Cucker and D. X. Zhou. *Learning Theory: An Approximation Theory View Point*. Cambridge University Press, 2007.
- P. Diaconis. Bayesian numerical analysis. *Statistical decision theory and related topics IV*, 1: 163–175, 1988.
- J. Dick, F. Y. Kuo, and I. H. Sloan. High dimensional numerical integration - the Quasi-Monte Carlo way. *Acta Numerica*, 22(133-288), 2013.
- M. F. Driscoll. The reproducing kernel Hilbert space structure of the sample paths of a gaussian process. *Zeitschrift für Wahrscheinlichkeitstheorie und verwandte Gebiete*, 26(4): 309–316, 1973.
- R. M. Dudley. *Real Analysis and Probability*. Cambridge University Press, 2002.
- D. E. Edmunds and H. Triebel. *Function Spaces, Entropy Numbers, Differential Operators*. Cambridge University Press, Cambridge, 1996.
- S. Flaxman, D. Sejdinovic, J. Cunningham, and S. Filippi. Bayesian learning of kernel embeddings. In *Uncertainty in Artificial Intelligence (UAI)*, pages 182–191, 2016.
- K. Fukumizu, F. Bach, and M. Jordan. Dimensionality reduction for supervised learning with reproducing kernel hilbert spaces. *Journal of Machine Learning Research*, 5:73–99, 2004.
- K. Fukumizu, A. Gretton, X. Sun, and B. Schölkopf. Kernel measures of conditional dependence. In *Advances in Neural Information Processing Systems 20*, pages 489–496, 2008.
- J. M. Gonzalez-Barrios and R. M. Dudley. Metric entropy conditions for an operator to be of trace class. *Proceedings of the American Mathematical Society*, 118(1):175–180, 1993.

- A. Gretton, O. Bousquet, A. Smola, and B. Schölkopf. Measuring statistical dependence with Hilbert-Schmidt norms. In S. Jain, H. U. Simon, and E. Tomita, editors, *Algorithmic Learning Theory*, volume 3734 of *Lecture Notes in Computer Science*, pages 63–77, Berlin/Heidelberg, 2005. Springer-Verlag.
- A. Gretton, K. Fukumizu, C. Teo, L. Song, B. Schoelkopf, and A. Smola. A kernel statistical test of independence. In *Advances in Neural Information Processing Systems 20*, pages 585–592, 2008.
- A. Gretton, K. Borgwardt, M. Rasch, B. Schölkopf, and A. Smola. A kernel two-sample test. *Journal of Machine Learning Research*, 13:723–773, 2012.
- T. Gunter, M. A. Osborne, R. Garnett, P. Hennig, and S. J. Roberts. Sampling for inference in probabilistic models with fast Bayesian quadrature. In Z. Ghahramani, M. Welling, C. Cortes, N. D. Lawrence, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems 27*, pages 2789–2797. Curran Associates, Inc., 2014.
- P. Hennig, M. A. Osborne, and M. Girolami. Probabilistic numerics and uncertainty in computations. *Proceedings of the Royal Society of London A: Mathematical, Physical and Engineering Sciences*, 471(2179), 2015.
- F. J. Hickernell. A generalized discrepancy and quadrature error bound. *Mathematics of Computation of the American Mathematical Society*, 67(221):299–322, 1998.
- T. Hofmann, B. Schölkopf, and A. J. Smola. Kernel methods in machine learning. *Annals of Statistics*, 36(3):1171–1220, 2008.
- F. Huszár and D. Duvenaud. Optimally-weighted herding is Bayesian quadrature. In *Uncertainty in Artificial Intelligence*, pages 377–385, 2012.
- S. Janson. *Gaussian Hilbert Spaces*. Cambridge University Press, 1997.
- M. Kanagawa, B. K. Sriperumbudur, and K. Fukumizu. Convergence guarantees for kernel-based quadrature rules in misspecified settings. In D. D. Lee, M. Sugiyama, U. V. Luxburg, I. Guyon, and R. Garnett, editors, *Advances in Neural Information Processing Systems 29*, pages 3288–3296. Curran Associates, Inc., 2016.
- M. Kanagawa, B. K. Sriperumbudur, and K. Fukumizu. Convergence analysis of deterministic kernel-based quadrature rules in misspecified settings. *Arxiv e-prints*, arXiv:1709.00147v1 [math.NA], Sept. 2017.
- T. Karvonen, C. J. Oates, and S. Särkkä. A Bayes-Sard cubature method. *Arxiv e-prints*, arXiv:1804.03016v3 [stat.ME], 2018.
- G. S. Kimeldorf and G. Wahba. A correspondence between Bayesian estimation on stochastic processes and smoothing by splines. *The Annals of Mathematical Statistics*, 41(2):495–502, 1970.
- A. N. Kolmogorov. Interpolation and extrapolation of stationary random sequences. math. 5. *Bull. Moscow Univ., Moscow*, 1941.

- F. M. Larkin. Gaussian measure in Hilbert space and applications in numerical analysis. *Rocky Mountain Journal of Mathematics*, 2(3):379–422, 1972.
- H. C. L. Law, D. Sutherland, D. Sejdinovic, and S. Flaxman. Bayesian approaches to distribution regression. In A. Storkey and F. Perez-Cruz, editors, *Proceedings of the Twenty-First International Conference on Artificial Intelligence and Statistics*, volume 84 of *Proceedings of Machine Learning Research*, pages 1167–1176. PMLR, 2018.
- F. Lindgren, H. Rue, and J. Lindsröm. An explicit link between Gaussian fields and Gaussian Markov random fields: the stochastic partial differential equation approach. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 73(4):423–498, 2011.
- Q. Liu and J. Lee. Black-box importance sampling. In A. Singh and J. Zhu, editors, *Proceedings of the 20th International Conference on Artificial Intelligence and Statistics*, volume 54 of *Proceedings of Machine Learning Research*, pages 952–961. PMLR, 2017.
- Q. Liu, J. Lee, and M. I. Jordan. A kernelized Stein discrepancy for goodness-of-fit tests. In M. F. Balcan and K. Q. Weinberger, editors, *Proceedings of The 33rd International Conference on Machine Learning*, volume 48 of *Proceedings of Machine Learning Research*, pages 276–284. PMLR, 2016.
- M. N. Lukić and J. H. Beder. Stochastic processes with sample paths in reproducing kernel Hilbert spaces. *Transactions of the American Mathematical Society*, 353(10):729–748, 2001.
- B. Matérn. Spatial variation. *Meddelanden fran Statens Skogsforskningsinstitut*, 49(5), 1960.
- G. Matheron. *Traité de géostatistique appliquée. 1 (1962)*, volume 1. Editions Technip, 1962.
- M. Meister and I. Steinwart. Optimal learning rates for localized SVMs. *Journal of Machine Learning Research*, 17(194):1–44, 2016.
- J. Mercer. Functions of positive and negative type, and their connection with the theory of integral equations. *Philosophical Transactions of the Royal Society A*, 209:415–446, 1909.
- K. Muandet, B. Sriperumbudur, K. Fukumizu, A. Gretton, and B. Schölkopf. Kernel mean shrinkage estimators. *Journal of Machine Learning Research*, 17(48):1–41, 2016.
- K. Muandet, K. Fukumizu, B. K. Sriperumbudur, and B. Schölkopf. Kernel mean embedding of distributions : A review and beyond. *Foundations and Trends in Machine Learning*, 10(1–2):1–141, 2017.
- R. M. Neal. Regression and classification using Gaussian process priors. In e. a. J. M. Bernardo, editor, *Bayesian Statistics 6*, pages 475–501. Oxford University Press, 1998.
- E. Novak and H. Wóznickowski. *Tractability of Multivariate Problems, Vol. I: Linear Information*. EMS, 2008.
- E. Novak and H. Wóznickowski. *Tractability of Multivariate Problems, Vol. II: Standard Information for Functionals*. EMS, 2010.
- C. J. Oates, M. Girolami, and N. Chopin. Control functionals for Monte Carlo integration. *Journal of the Royal Statistical Society, Series B*, 79(2):323–380, 2017.

- A. O’Hagan. Bayes–Hermite quadrature. *Journal of Statistical Planning and Inference*, 29(3): 245–260, 1991.
- E. Parzen. An approach to time series analysis. *The Annals of Mathematical Statistics*, 32(4): 951–989, 1961.
- N. Pfister, P. Bühlmann, B. Schölkopf, and J. Peters. Kernel-based tests for joint independence. *Journal of the Royal Statistical Society, Series B*, 80(1):5–31, 2017.
- C. Rasmussen and C. Williams. *Gaussian Processes for Machine Learning*. MIT Press, 2006.
- K. Ritter. *Average-Case Analysis of Numerical Problems*. Springer, 2000.
- R. Schaback and H. Wendland. Kernel techniques: From machine learning to meshless methods. *Acta Numerica*, 15:543–639, 2006.
- M. Scheuere, R. Schaback, and M. Schlather. Interpolation of spatial data - A stochastic or a deterministic problem? *European Journal of Applied Mathematics*, 24:601–629, 2013.
- B. Schölkopf and A. J. Smola. *Learning with Kernels: Support Vector Machines, Regularization, Optimization, and Beyond*. MIT press, 2002.
- B. Schölkopf, R. Herbrich, and A. Smola. A generalized representer theorem. In *Computational Learning Theory. COLT 2001. Lecture Notes in Computer Science*, volume 2111, pages 416–426. Springer, 2001.
- B. Schölkopf, K. Tsuda, and J. P. Vert. *Kernel Methods in Computational Biology*. MIT Press, 2004.
- D. Sejdinovic, B. Sriperumbudur, A. Gretton, and K. Fukumizu. Equivalence of distance-based and RKHS-based statistics in hypothesis testing. *Annals of Statistics*, 41(5):2263–2702, 2013.
- A. Smola, A. Gretton, L. Song, and B. Schölkopf. A Hilbert space embedding for distributions. In *Proceedings of the International Conference on Algorithmic Learning Theory*, volume 4754, pages 13–31. Springer, 2007.
- L. Song, A. Smola, A. Gretton, J. Bedo, and K. Borgwardt. Feature selection via dependence maximization. *Journal of Machine Learning Research*, 13:1393–1434, 2012.
- B. K. Sriperumbudur, A. Gretton, K. Fukumizu, B. Schölkopf, and G. R. Lanckriet. Hilbert space embeddings and metrics on probability measures. *Journal of Machine Learning Research*, 11:1517–1561, 2010.
- B. K. Sriperumbudur, K. Fukumizu, and G. R. Lanckriet. Universality, characteristic kernels and RKHS embedding of measures. *Journal of Machine Learning Research*, 12:2389–2410, 2011.
- C. Stein. Inadmissibility of the usual estimator for the mean of a multivariate normal distribution. In *Proceedings of the Third Berkeley Symposium on Mathematical Statistics and Probability*, volume 1, pages 197–206. University of California Press, 1956.

- E. M. Stein. *Singular Integrals and Differentiability Properties of Functions*. Princeton University Press, Princeton, NJ, 1970.
- M. L. Stein. *Interpolation of Spatial Data*. Springer-Verlag, New York, 1999.
- I. Steinwart. Convergence types and rates in generic Karhunen-Loéve expansions with applications to sample path properties. *ArXiv e-prints*, arXiv:1403.1040v3 [math.PR], Mar. 2017.
- I. Steinwart and A. Christmann. *Support Vector Machines*. Springer, 2008.
- I. Steinwart and C. Scovel. Mercer’s theorem on general domains: on the interaction between measures, kernels, and RKHS. *Constructive Approximation*, 35:363–417, 2012.
- I. Steinwart, D. Hush, and C. Scovel. Optimal rates for regularized least squares regression. In S. Dasgupta and A. Klivans, editors, *Proceedings of the 22nd Annual Conference on Learning Theory*, pages 79–93, 2009.
- C. J. Stone. Optimal rates of convergence for nonparametric estimators. *The Annals of Statistics*, 8(6):1348–1360, 1980.
- A. M. Stuart. Inverse problems: A Bayesian perspective. *Acta Numerica*, 19:451–559, 2010.
- A. M. Stuart and A. L. Teckentrup. Posterior consistency for Gaussian process approximations of Bayesian posterior distributions. *Mathematics of Computation*, 87:721–753, 2018.
- Z. Szabó and B. K. Sriperumbudur. Characteristic and universal tensor product kernels. *ArXiv e-prints*, arXiv:1708.08157v3 [stat.ML], May 2018.
- G. J. Székely and M. L. Rizzo. Brownian distance covariance. *The Annals of Applied Statistics*, 3(4):1236–1265, 2009.
- I. Tolstikhin, B. K. Sriperumbudur, and K. Muandet. Minimax estimation of kernel mean embeddings. *Journal of Machine Learning Research*, 18(86):1–47, 2017.
- H. Triebel. *Theory of Function Spaces III*. Birkhäuser Verlag, 2006.
- A. B. Tsybakov. *Introduction to Nonparametric Estimation*. Springer, 2008.
- R. Tuo and C. F. J. Wu. A theoretical framework for calibration in computer models: Parametrization, estimation and convergence properties. *SIAM/ASA Journal on Uncertainty Quantification*, 4(1):739–766, 2016.
- A. van der Vaart and H. van Zanten. Information rates of nonparametric Gaussian process methods. *Journal of Machine Learning Research*, 12:2095–2119, 2011.
- A. W. van der Vaart and J. H. van Zanten. Reproducing kernel Hilbert spaces of Gaussian priors. *IMS Collections, Pushing the Limits of Contemporary Statistics: Contributions in Honor of Jayanta K. Ghosh*, 3:200–222, 2008.
- G. Wahba. *Spline Models for Observational Data*. Number 59 in CBMS-NSF Regional Conferences series in applied mathematics. SIAM, 1990.

- H. Wendland. *Scattered Data Approximation*. Cambridge University Press, Cambridge, UK, 2005.
- H. Wendland and C. Rieger. Approximate interpolation with applications to selecting smoothing parameters. *Numerische Mathematik*, 101(4):729–748, 2005.
- P. Whittle. On stationary processes in the plane. *Biometrika*, 41(3/4):434–449, 1954.
- Z. Wu and R. Schaback. Local error estimates for radial basis function interpolation of scattered data. *IMA journal of Numerical Analysis*, 13(1):13–27, 1993.
- M. Yamada, W. Jitkrittum, L. Sigal, E. P. Xing, and M. Sugiyama. High-dimensional feature selection by feature-wise non-linear lasso. *Neural Computation*, 26(1):185–207, 2014.
- M. Yamada, Y. Umezū, K. Fukumizu, and I. Takeuchi. Post selection inference with kernels. In A. Storkey and F. Perez-Cruz, editors, *Proceedings of the Twenty-First International Conference on Artificial Intelligence and Statistics*, volume 84 of *Proceedings of Machine Learning Research*, pages 152–160. PMLR, 2018.
- Q. Zhang, S. Filippi, A. Gretton, and D. Sejdinovic. Large-scale kernel methods for independence testing. *Statistics and Computing*, 28(1):113–130, 2018.