

# Parameterization Schemes and Their Quality in Kernel Interpolation

Michael McCourt

January 2, 2014

## Abstract

When performing interpolation with a kernel basis, the choice of kernel can play a significant role in the accuracy of the interpolation. Unfortunately, choosing a “good” kernel is a difficult proposition because of the wide variety of kernels available. To simplify this process, often a family of kernels is considered which differ by one or more free parameters while still retaining many of the same interpolation properties. Gaussians are an example of one such kernel, whose shape parameter allows for variable localization. Our goal is to consider several parameterization schemes for Gaussians and to develop a mechanism for comparing these schemes so that we can determine which is effective under what circumstances.

## 1 Introduction

Kernel interpolation is discussed in [2, 8, 7] and many other books. It is a useful tool for approximating scattered data and appears in both approximation theory (through reproducing kernel Hilbert spaces) and in spatial statistics (through Kriging). We believe that by leveraging knowledge from both of these fields we can improve upon the existing literature. There is another writeup of this content by Fred Hickernell available at

<http://math.ucdenver.edu/~mmccourt/hsp.pdf>

although the notation is somewhat in conflict with what I have written here. Still, it is a fantastic writeup and may help to supplement this discussion.

In both these settings, we are provided data

$$\mathcal{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}, \quad \mathbf{y} = \begin{pmatrix} y_1 \\ \vdots \\ y_N \end{pmatrix}, \quad \mathbf{x}_i \neq \mathbf{x}_j \text{ for } i \neq j,$$

where  $\mathbf{x}_i \in \Omega \subseteq \mathbb{R}^d$  and  $y_i \in \mathbb{R}$ . Our goal is to somehow make “good” guesses about values of  $y$  at  $\mathbf{x} \in \Omega$  locations where data does not exist.

A kernel  $K : \Omega \times \Omega \rightarrow \mathbb{R}$  is any symmetric function of two variables. For this work, we restrict  $K$  to the set of *positive definite* kernels, which are kernels such that

$$\int_{\Omega} \int_{\Omega} K(\mathbf{x}, \mathbf{z}) v(\mathbf{x}) v(\mathbf{z}) d\mathbf{x} d\mathbf{z} > 0$$

for any  $v \in L^2(\Omega)$  not identically zero. This has another (more useful) definition involving the eigenvalues of a specific linear operator, but we’ll worry about that if we need it.

Kernels can take lots of different forms and have lots of different properties, but probably the most common form of a kernel is a *radial basis function*, which is a kernel of the form

$$K(\mathbf{x}, \mathbf{z}) = \phi(\|\mathbf{x} - \mathbf{z}\|), \quad \mathbf{x}, \mathbf{z} \in \Omega.$$

Notice that  $\phi : \mathbb{R}^+ \rightarrow \mathbb{R}$  is a function of one variable instead of  $\mathbb{R}^d$  which is one of the many reasons why radial kernels are nice. The kernel of main interest right now will be the Gaussian kernel

$$K(\mathbf{x}, \mathbf{z}) = \exp(-\varepsilon^2 \|\mathbf{x} - \mathbf{z}\|^2); \tag{1.1}$$

the norm used above is the 2-norm and the value  $\varepsilon > 0$  is a *shape parameter* which can be chosen freely. The “correct” choice of  $\varepsilon$  is an important component in making “good” guesses about unobserved  $(\mathbf{x}, y)$  values.

## 1.1 Scattered Data Interpolation

In the interpolation setting, we assume that data we are given  $\{(\mathbf{x}_i, y_i)\}_{i=1}^N$  has been generated by some function  $f$  such that  $y_i = f(\mathbf{x}_i)$ . That function is deterministic, and we assume  $f \in L^2(\Omega)$ ; it is also common to assume that  $f$  is in a certain Hilbert space, but we will state when that assumption is made explicitly.

To make guesses of  $y$  values at unobserved  $\mathbf{x}$  locations, we create an interpolant  $s$  of  $f$ , and then evaluate  $s(\mathbf{x})$ . Because we are performing kernel interpolation,  $s$  uses a linear combination of kernel functions

$$s(\mathbf{x}) = \sum_{k=1}^{\hat{N}} c_k K(\mathbf{x}, \mathbf{z}_k),$$

where  $\mathbf{z}_k$  for  $1 \leq k \leq \hat{N}$  are the so-called *kernel centers* that distinguish the different elements of the basis.

Often times, we will prefer to write the interpolant as a vector inner product,

$$s(\mathbf{x}) = \mathbf{k}(\mathbf{x})^T \mathbf{c}, \tag{1.2}$$

using

$$\mathbf{k}(\mathbf{x}) = \begin{pmatrix} K(\mathbf{x}, \mathbf{x}_1) \\ \vdots \\ K(\mathbf{x}, \mathbf{x}_N) \end{pmatrix}, \quad \mathbf{c} = \begin{pmatrix} c_1 \\ \vdots \\ c_N \end{pmatrix}.$$

To be explicit, we should probably write  $\mathbf{k}_{\mathcal{X}}(\mathbf{x})$  to remind ourselves that our basis is data-dependent, but we will omit it for now. To solve the interpolation problem we require  $s(\mathbf{x}_i) = y_i$  for  $1 \leq i \leq N$ , which produces the linear system

$$\begin{pmatrix} y_1 \\ \vdots \\ y_N \end{pmatrix} = \begin{pmatrix} s(\mathbf{x}_1) \\ \vdots \\ s(\mathbf{x}_N) \end{pmatrix} = \begin{pmatrix} \mathbf{k}(\mathbf{x}_1)^T \mathbf{c} \\ \vdots \\ \mathbf{k}(\mathbf{x}_N)^T \mathbf{c} \end{pmatrix} = \begin{pmatrix} \mathbf{k}(\mathbf{x}_1)^T \\ \vdots \\ \mathbf{k}(\mathbf{x}_N)^T \end{pmatrix} \mathbf{c},$$

or, more succinctly,

$$\mathbf{K} \mathbf{c} = \mathbf{y}, \quad \mathbf{K} = \begin{pmatrix} \mathbf{k}(\mathbf{x}_1)^T \\ \vdots \\ \mathbf{k}(\mathbf{x}_N)^T \end{pmatrix}. \tag{1.3}$$

Because  $K$  is a positive definite kernel, we are guaranteed that  $\mathbf{K}$  is a symmetric positive definite matrix. This implies that kernel-interpolation is well-defined in any dimension, in contrast to polynomial interpolation. Using this definition of  $\mathbf{c}$  in (1.2) gives

$$s(\mathbf{x}) = \mathbf{k}(\mathbf{x})^T \mathbf{K}^{-1} \mathbf{y}, \tag{1.4}$$

which explicitly shows how values at unobserved locations are a weighted averaging of the given data.

To try to estimate the error associated with your kernel interpolation, some theory exists to help, although it may not be as useful as we would like. We can bound the pointwise error of the interpolant (a big deal) by the native space norm of the function (incomputable):

$$|f(\mathbf{x}) - s(\mathbf{x})| \leq P_{K, \mathcal{X}}(\mathbf{x}) \|f\|_{\mathcal{H}_K(\Omega)}, \quad f \in \mathcal{H}_K(\Omega), \tag{1.5}$$

where  $P_{K, \mathcal{X}}(\mathbf{x})$  is the so-called *power function*. The norm in  $\mathcal{H}_K(\Omega)$  is called the *native space* norm, and that means the Hilbert space norm of  $f$  for the Hilbert space induced by the kernel  $K$ . Computing that is not possible without the Hilbert-Schmidt decomposition of the kernel, but note that the value of  $\|f\|_{\mathcal{H}_K(\Omega)}$  is dependent on the kernel  $K$ . The power function is defined as

$$P_{K, \mathcal{X}}(\mathbf{x}) = \sqrt{K(\mathbf{x}, \mathbf{x}) - \mathbf{k}(\mathbf{x})^T \mathbf{K}^{-1} \mathbf{k}(\mathbf{x})} \tag{1.6}$$

and comes about from some standard reproducing kernel manipulations [2] written up in the appendix.

In the definition of the Gaussian (1.1), there is an  $\varepsilon$  value which serves as the shape parameter of the basis: small  $\varepsilon$  produces flat kernels, and large  $\varepsilon$  produces peaked kernels. These differences can be seen in Figure 1.

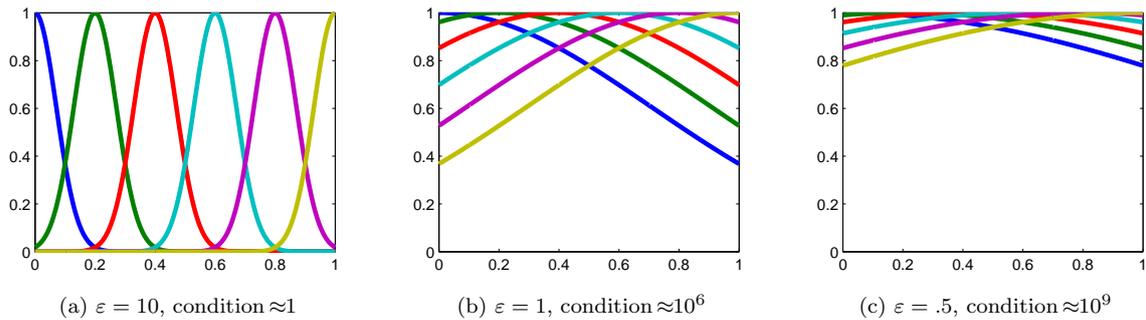


Figure 1: Gaussian basis functions centered at 6 points in  $[0,1]$ . The choice of shape parameter affects the “width” or “locality” of the kernels. The condition of the  $K$  interpolation matrix is also listed.

### 1.1.1 Hilbert-Schmidt SVD

The following derivation took place in [1] so I’m not going to go too in depth here. I just want to introduce the notation for use later as necessary.

Positive definite kernels have eigenvalues and eigenfunctions defined by the Hilbert-Schmidt integral operator

$$\int_{\Omega} K(\mathbf{x}, \mathbf{z}) \varphi_n(\mathbf{z}) \rho(\mathbf{z}) d\mathbf{z} = \lambda_n \varphi_n(\mathbf{x}), \quad (1.7)$$

where  $\rho$  is a suitable weight function. For Gaussians (our main focus right now) this discussion appears primarily in [3], where the values of  $\rho$ ,  $\lambda_n$  and  $\varphi_n$  are explicitly stated. We assume that  $\Omega = \mathbb{R}^d$ . It is relevant to note that as  $\varepsilon \rightarrow 0$ ,  $\lambda_n \approx \varepsilon^{2n}$ , which is the main cause of Gaussian ill-conditioning.

These eigenvalues are all positive, because  $K$  is a positive definite kernel, and decreasing, i.e.,  $\lambda_1 \geq \lambda_2 \geq \dots > 0$ ; their (infinite) sum is finite because the Hilbert-Schmidt operator is a trace class operator. The eigenfunctions satisfy the orthonormality property

$$\int_{\Omega} \varphi_m(\mathbf{x}) \varphi_n(\mathbf{x}) \rho(\mathbf{x}) d\mathbf{x} = \begin{cases} 1 & m = n, \\ 0 & m \neq n, \end{cases}$$

which might also be written as  $\langle \varphi_m, \varphi_n \rangle_{L_2(\rho)} = \delta_{m,n}$ . Another important orthogonality property of the eigenfunctions involves the Hilbert space inner product. Recall first the reproducing property

$$\langle f, K(\cdot, \mathbf{x}) \rangle_{\mathcal{H}_K} = f(\mathbf{x}), \quad f \in \mathcal{H}_K$$

which was discussed and used in the Appendix. Because  $\varphi_m \in \mathcal{H}_K$  for  $n = 1, 2, \dots$ , and using (1.7) to write  $\varphi_m = \frac{1}{\lambda_m} \int_{\Omega} K(\cdot, \mathbf{z}) \varphi_m(\mathbf{z}) \rho(\mathbf{z}) d\mathbf{z}$ , we know that

$$\begin{aligned} \langle \varphi_m, \varphi_n \rangle_{\mathcal{H}_K} &= \left\langle \frac{1}{\lambda_m} \int_{\Omega} K(\cdot, \mathbf{z}) \varphi_m(\mathbf{z}) \rho(\mathbf{z}) d\mathbf{z}, \varphi_n \right\rangle_{\mathcal{H}_K} \\ &= \frac{1}{\lambda_m} \int_{\Omega} \langle K(\cdot, \mathbf{z}), \varphi_n \rangle_{\mathcal{H}_K} \varphi_m(\mathbf{z}) \rho(\mathbf{z}) d\mathbf{z} \\ &= \frac{1}{\lambda_m} \int_{\Omega} \varphi_m(\mathbf{z}) \varphi_n(\mathbf{z}) \rho(\mathbf{z}) d\mathbf{z} = \frac{\delta_{mn}}{\lambda_m} \end{aligned} \quad (1.8)$$

These eigenvalues/functions allow us to write the kernel  $K$  as a Mercer’s series (or Hilbert-Schmidt series)

$$K(\mathbf{x}, \mathbf{z}) = \sum_{m=1}^{\infty} \lambda_m \varphi_m(\mathbf{x}) \varphi_m(\mathbf{z})$$

which means that the kernel basis vector  $\mathbf{k}$  can be written as

$$\mathbf{k}(\mathbf{x})^T = \boldsymbol{\phi}(\mathbf{x})^T \boldsymbol{\Lambda} \boldsymbol{\Phi}^T$$

for  $\phi(\mathbf{x})^T = (\varphi_1(\mathbf{x}) \ \cdots \ \varphi_N(\mathbf{x}) \ \cdots)$ , which is an infinite length vector, and

$$\Lambda = \begin{pmatrix} \lambda_1 & & | & \\ & \ddots & & | \\ \hline & & \lambda_N & | \\ & & & \ddots \end{pmatrix} = \begin{pmatrix} \Lambda_1 & \\ & \Lambda_2 \end{pmatrix}, \quad \Phi = \begin{pmatrix} \phi(\mathbf{x}_1)^T \\ \vdots \\ \phi(\mathbf{x}_N)^T \end{pmatrix} = (\Phi_1 \ \Phi_2)$$

where  $\Lambda_1, \Phi_1 \in \mathbb{R}^{N \times N}$  and  $\Lambda_2$  and  $\Phi_2$  are the (infinite-sized) rest of the matrices  $\Lambda$  and  $\Phi$  respectively. Eigenfunctions can be chosen so that  $\Phi_1^{-1}$  exists, and, because  $K$  is a positive definite kernel,  $\lambda_n > 0$  for  $1 \leq n < \infty$  thus we know that  $\Lambda_1^{-1}$  exists. This allows us to write

$$\Lambda \Phi^T = \begin{pmatrix} \Lambda_1 & \\ & \Lambda_2 \end{pmatrix} \begin{pmatrix} \Phi_1^T \\ \Phi_2^T \end{pmatrix} = \begin{pmatrix} \Lambda_1 \Phi_1^T \\ \Lambda_2 \Phi_2^T \Phi_1^{-T} \Lambda_1^{-1} \end{pmatrix} \Lambda_1 \Phi_1^T.$$

Using this in our  $\mathbf{k}$  definition above gives

$$\mathbf{k}(\mathbf{x})^T = \underbrace{\phi(\mathbf{x})^T}_{\psi(\mathbf{x})^T} \begin{pmatrix} \Lambda_1 \Phi_1^T \\ \Lambda_2 \Phi_2^T \Phi_1^{-T} \Lambda_1^{-1} \end{pmatrix} \Lambda_1 \Phi_1^T = \psi(\mathbf{x})^T \Lambda_1 \Phi_1^T. \quad (1.9)$$

Substituting this into our interpolation matrix  $\mathbf{K}$  from (1.3) gives

$$\mathbf{K} = \begin{pmatrix} \mathbf{k}(\mathbf{x}_1)^T \\ \vdots \\ \mathbf{k}(\mathbf{x}_N)^T \end{pmatrix} = \begin{pmatrix} \psi(\mathbf{x}_1)^T \Lambda_1 \Phi_1^T \\ \vdots \\ \psi(\mathbf{x}_N)^T \Lambda_1 \Phi_1^T \end{pmatrix} = \begin{pmatrix} \psi(\mathbf{x}_1)^T \\ \vdots \\ \psi(\mathbf{x}_N)^T \end{pmatrix} \Lambda_1 \Phi_1^T = \Psi \Lambda_1 \Phi_1^T, \quad (1.10)$$

which is the *Hilbert-Schmidt SVD*. This is a useful decomposition because much of the ill-conditioning of  $\mathbf{K}$  resides in  $\Lambda_1$  (subject to some appropriate design choices discussed in [3]). Using (1.10) and (1.9) in (1.4) gives

$$s(\mathbf{x}) = \mathbf{k}(\mathbf{x})^T \mathbf{K}^{-1} \mathbf{y} = \psi(\mathbf{x})^T \Lambda_1 \Phi_1^T (\Psi \Lambda_1 \Phi_1^T)^{-1} \mathbf{y} = \psi(\mathbf{x})^T \Psi^{-1} \mathbf{y} \quad (1.11)$$

which shows that the Hilbert-Schmidt SVD allows for a change of basis from the unstable basis  $\mathbf{k}(\mathbf{x})$  to the stable basis  $\psi(\mathbf{x})$ .

## 1.2 Kriging

This section I am less confident about, especially regarding notation, so please correct as needed. We'll start with the idea of a probability space  $(\mathcal{W}, \mathcal{A}, P)$  where  $\mathcal{W}$  is the sample space of all possible outcomes,  $\mathcal{A}$  is a  $\sigma$ -algebra and  $P$  is a probability measure. Normally, I think statisticians and probabilists would use  $\Omega$  instead of  $\mathcal{W}$ , but we're already using  $\Omega$  for something else.

We need to define a parameter space  $\Omega$ , which for our purposes right now will be  $\Omega = \mathbb{R}^d$ . In numerical analysis, this is the domain over which we would consider evaluating our interpolant. A function  $Y : \Omega, \mathcal{W} \rightarrow \mathbb{R}$  (evaluated as  $Y(\mathbf{x}, \omega)$  for  $\mathbf{x} \in \Omega$  and  $\omega \in \mathcal{W}$ ) is a random field if, for every  $\mathbf{x} \in \Omega$ ,  $Y$  is an  $\mathcal{A}$ -measurable function of  $\omega$ . Our notation for this is

$$\text{Random Field: } Y = \{Y_{\mathbf{x}}\}_{\mathbf{x} \in \Omega}.$$

Note that

- For a fixed  $\mathbf{x}$ ,  $Y_{\mathbf{x}} = Y(\mathbf{x}, \cdot)$  is a random variable.
- For a fixed  $\omega$ ,  $y(\cdot) = Y(\cdot, \omega)$  is a deterministic function of  $\mathbf{x}$  referred to as a realization of the random field.

The mean of  $Y$  is a function  $\mu_Y$  which is defined at any point  $\mathbf{x} \in \Omega$  as

$$\mu_Y(\mathbf{x}) = \mathbb{E}(Y_{\mathbf{x}}) = \int_{\mathcal{W}} Y_{\mathbf{x}}(\omega) dP(\omega) = \int_{\mathbb{R}} y dF_{Y_{\mathbf{x}}}(y)$$

where  $F_{Y_{\mathbf{x}}}$  is the cumulative distribution function of  $Y_{\mathbf{x}}$ . For our purposes we will assume that  $Y_{\mathbf{x}}$  is continuous, thus we can write

$$\mu_Y(\mathbf{x}) = \int_{\mathbb{R}} yp_{Y_{\mathbf{x}}}(y)dy$$

for density function  $p_{Y_{\mathbf{x}}}$ . Likewise, we will define the covariance kernel of  $Y$  as

$$K(\mathbf{x}, \mathbf{z}) = \text{Cov}(Y_{\mathbf{x}}, Y_{\mathbf{z}}) = \mathbb{E}((Y_{\mathbf{x}} - \mu_Y(\mathbf{x}))(Y_{\mathbf{z}} - \mu_Y(\mathbf{z}))) = \mathbb{E}(Y_{\mathbf{x}}Y_{\mathbf{z}}) - \mu_Y(\mathbf{x})\mu_Y(\mathbf{z}), \quad (1.12)$$

after some appropriate manipulations. We are not using the notation  $K_Y$  here, unlike  $\mu_Y$ , to later emphasize the similarities between Kriging and kernel approximation.

Our assumption here is that the random field  $Y$  is a *Gaussian random field*, which we denote as

$$Y \sim GF(\mu_Y, K).$$

This implies that, for a finite set of points  $\mathcal{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\} \in \Omega$ , the vector of random variables  $\mathbf{Y}_{\mathcal{X}} = (Y_{\mathbf{x}_1} \ \dots \ Y_{\mathbf{x}_N})^T$  has the distribution

$$\mathbf{Y}_{\mathcal{X}} \sim \mathcal{N}(\boldsymbol{\mu}, \mathbf{K}), \quad (1.13)$$

where  $\boldsymbol{\mu} = \mathbb{E}(\mathbf{Y}_{\mathcal{X}})$  and  $(\mathbf{K})_{i,j} = \text{Cov}(Y_{\mathbf{x}_i}, Y_{\mathbf{x}_j})$ . This multivariate normal distribution has the density

$$p_{\mathbf{Y}_{\mathcal{X}}}(\mathbf{y}) = \frac{1}{\sqrt{2\pi \det(\mathbf{K})}} \exp\left(-\frac{1}{2}(\mathbf{y} - \boldsymbol{\mu})^T \mathbf{K}^{-1}(\mathbf{y} - \boldsymbol{\mu})\right), \quad (1.14)$$

and the term  $(\mathbf{y} - \boldsymbol{\mu})^T \mathbf{K}^{-1}(\mathbf{y} - \boldsymbol{\mu})$  is sometimes referred to as the Mahalanobis distance. Note that  $\mathbf{K}$  will be symmetric positive definite for a positive definite covariance kernel  $K$ , which must be positive definite or else a negative variance could occur. As a result,  $\mathbf{K}^{-1}$  must exist, though it may be ill-conditioned.

At this point, we are going to restrict our concern to **zero-mean** Gaussian processes, which demands that  $\mu_Y \equiv 0$ ; in Kriging literature, this is referred to as “simple Kriging”. [\[MJM\] What penalties does this restriction incur?](#) This restriction simplifies the situation significantly, most notably by eliminating the mean terms from (1.12),

$$K(\mathbf{x}, \mathbf{z}) = \text{Cov}(Y_{\mathbf{x}}, Y_{\mathbf{z}}) = \mathbb{E}(Y_{\mathbf{x}}Y_{\mathbf{z}}), \quad (1.15)$$

and from (1.14)

$$p_{\mathbf{Y}_{\mathcal{X}}}(\mathbf{y}) = \frac{1}{\sqrt{2\pi \det(\mathbf{K})}} \exp\left(-\frac{1}{2}\mathbf{y}^T \mathbf{K}^{-1}\mathbf{y}\right), \quad (1.16)$$

This also produces a duality between the reproducing kernel Hilbert space  $\mathcal{H}_K(\Omega)$  defined in Section 1.1 and the Hilbert space  $\mathcal{H}_Y$  which is the set of all linear combinations of random variables  $Y_{\mathbf{x}}$  together with their  $L_2(\Omega, \mathcal{A}, P)$ -limits. The *Loève representation theorem* says that the inner products in these Hilbert spaces are identical, i.e.,

$$\langle Y_{\mathbf{x}}, Y_{\mathbf{z}} \rangle_{\mathcal{H}_Y} = \mathbb{E}(Y_{\mathbf{x}}Y_{\mathbf{z}}) = K(\mathbf{x}, \mathbf{z}) = \langle K(\cdot, \mathbf{x}), K(\cdot, \mathbf{z}) \rangle_{\mathcal{H}_K},$$

where the final identity occurs because of the reproducing property of  $K$ . We believe that this implies a connection between the eigenfunctions  $\varphi_m$  defined in Section 1.1.1 and elements of the Karhunen-Loève expansion, but I’m not going to write about that here.

The biggest advantage of assuming zero-mean Gaussian fields comes when making predictions for unobserved values of  $Y_{\mathbf{x}_0}$ ,  $\mathbf{x}_0 \notin \mathcal{X}$ . After a lengthy derivation (which can be found in a talk I gave at one point) involving the joint distribution of  $(\mathbf{Y}_{\mathcal{X}}, Y_{\mathbf{x}_0})$ , it can be shown that the conditional density for the random variable  $Y_{\mathbf{x}_0}$  given a realization  $\mathbf{y}$  of the multivariate normal random variable  $\mathbf{Y}_{\mathcal{X}}$  is

$$p_{Y_{\mathbf{x}_0}}(y|\mathbf{Y}_{\mathcal{X}} = \mathbf{y}) \propto \exp\left(-\frac{1}{2}(y - \bar{\mu})C^{-1}(y - \bar{\mu})\right), \quad (1.17)$$

where  $\bar{\mu} = \mu_Y(\mathbf{x}_0) + \mathbf{k}(\mathbf{x})^T \mathbf{K}^{-1}(\mathbf{y} - \boldsymbol{\mu})$  and  $C = K(\mathbf{x}_0, \mathbf{x}_0) - \mathbf{k}(\mathbf{x}_0)^T \mathbf{K}^{-1} \mathbf{k}(\mathbf{x}_0)$ . This implies that  $Y_{\mathbf{x}_0}|\mathbf{Y}_{\mathcal{X}} = \mathbf{y}$  is a normal random variable (no surprise there, that the conditional of a normal is normal) with mean  $\bar{\mu}$  and variance  $C$ ; it is not a covariance here because we consider only one point  $\mathbf{x}_0$ .

By enforcing our zero-mean assumption, both  $\mu_Y(\mathbf{x}_0) = 0$  and  $\boldsymbol{\mu} = 0$ , meaning that

$$Y_{\mathbf{x}_0} | \mathbf{Y}_{\mathcal{X}} = \mathbf{y} \sim \mathcal{N} \left( \mathbf{k}(\mathbf{x}_0)^T \mathbf{K}^{-1} \mathbf{y}, K(\mathbf{x}_0, \mathbf{x}_0) - \mathbf{k}(\mathbf{x}_0)^T \mathbf{K}^{-1} \mathbf{k}(\mathbf{x}_0) \right). \quad (1.18)$$

Thus, the best linear unbiased predictor for the zero-mean Gaussian field  $Y$  at a point  $\mathbf{x}_0 \in \Omega$  is

$$\mathbb{E}(Y_{\mathbf{x}_0} | \mathbf{Y}_{\mathcal{X}} = \mathbf{y}) = \mathbf{k}(\mathbf{x}_0)^T \mathbf{K}^{-1} \mathbf{y}, \quad (1.19)$$

which through some miracle of mathematics is identical to (1.2) if you consider the function values in the numerical analysis setting  $\mathbf{y}$  to be equal to the Gaussian random field realization  $\mathbf{y}$ . In another awesome twist, the power function (1.6) evaluated at  $\mathbf{x}_0$  perfectly matches the so-called Kriging variance:

$$\text{Var}(Y_{\mathbf{x}_0} | \mathbf{Y}_{\mathcal{X}} = \mathbf{y}) = K(\mathbf{x}_0, \mathbf{x}_0) - \mathbf{k}(\mathbf{x}_0)^T \mathbf{K}^{-1} \mathbf{k}(\mathbf{x}_0) = P_{K, \mathcal{X}}(\mathbf{x}_0)^2.$$

This suggests that minimizing the power function, a wholly analytic approach to reducing the error in the interpolant, has the effect of minimizing the variance of our Kriging prediction, which is a wholly statistical approach to reducing the “error” in the prediction.

### 1.3 Parameterizing Kernels

So far, we have introduced two ways that we can approach this scattered data interpolation problem, which under certain circumstances yield the same result. We have assumed that the kernel which defines the interpolating basis  $\mathbf{k}(\mathbf{x})$ , and defines the covariance of the Gaussian random field  $Y$ , is a Gaussian kernel (1.1).

Within that Gaussian kernel  $K(\mathbf{x}, \mathbf{z}) = \exp(-\varepsilon^2 \|\mathbf{x} - \mathbf{z}\|^2)$  is a free parameter  $\varepsilon > 0$  which we have until now largely ignored. It can be proved that any value of  $\varepsilon$  will still produce a positive definite interpolation matrix  $\mathbf{K}$ , thus the interpolant (1.2) is guaranteed to exist and the covariance of  $Y_{\mathcal{X}}$  will be a positive definite matrix.

This theoretical result does not suggest, however, that all  $\varepsilon$  values will produce equally effective interpolants and predictions, nor does it suggest that the matrices  $\mathbf{K}$  will be well conditioned. Indeed, the predictive capacity of  $s(\mathbf{x}) = \mathbf{k}(\mathbf{x})^T \mathbf{K}^{-1} \mathbf{y}$  is greatly dependent on the shape parameter  $\varepsilon$ , as is indicated in Figure 2.

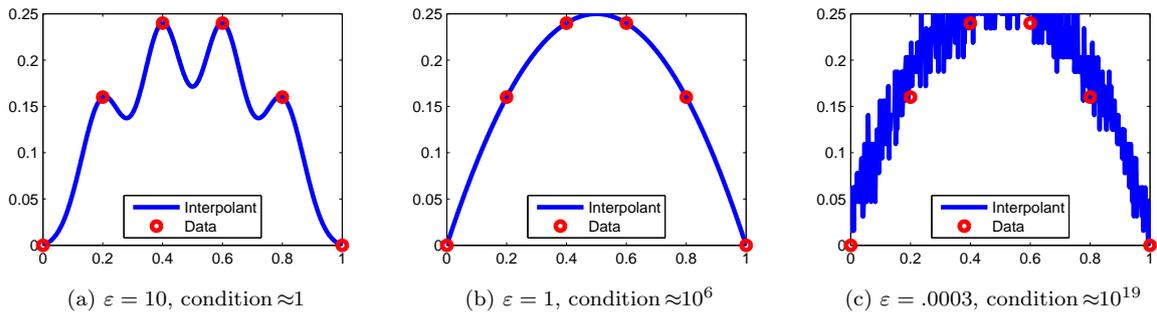


Figure 2: Different shape parameters will produce interpolants that look different. When  $\varepsilon$  is too large, the basis functions are too localized and predictions suffer. When  $\varepsilon$  is too small,  $\mathbf{K}$  is very ill-conditioned, yielding severe cancellation errors, and possibly preventing  $s(\mathbf{x}_i) = y_i$  for some of the given data. The function which created this data is  $f(x) = x(1 - x)$ .

The severe ill-conditioning which cripples the result in Figure 2c can be overcome with (1.11). But even with a stable result, as  $\varepsilon \rightarrow 0$  the Gaussian interpolant approaches the polynomial interpolant (or some polynomial approximation in higher dimensions) meaning that kernels have the potential to exceed polynomial accuracy for  $\varepsilon > 0$ . This is, at least in part, a result of the Runge phenomenon: polynomials are global basis functions, so because Gaussians have the ability to localize they have the ability to minimize the Runge phenomenon. This concept is demonstrated graphically in Figure 3.

Figure 3b demonstrates the opportunity for both success and failure when choosing  $\varepsilon$ . There is a minimum error attainable for  $\varepsilon \approx .4642$ , but if  $\varepsilon$  is chosen slightly smaller or larger the error could be orders of magnitude worse. Of course, this graph can only be created with the true function that generated the data, and that will be unavailable in many applications. This motivates the search for parameterization schemes that may allow us to find good  $\varepsilon$  values using just the available data.

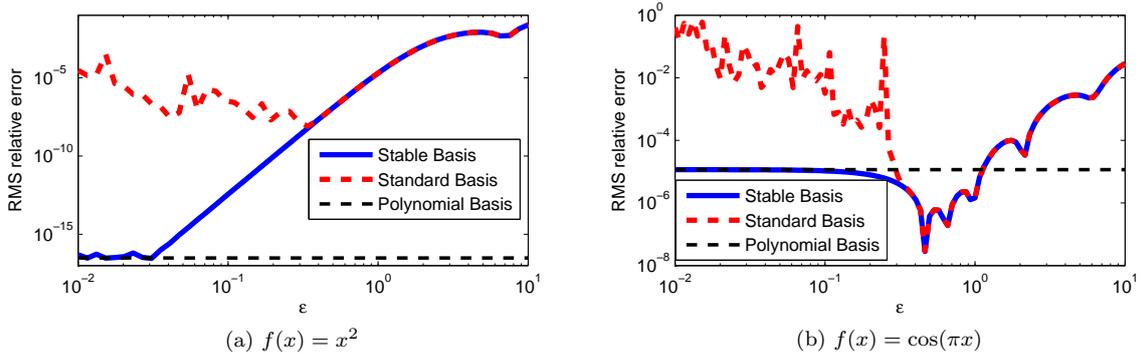


Figure 3: In this example, 10 evenly spaced points in  $(-1, 1)$  are sampled and Gaussian interpolation is conducted with a range of  $\varepsilon$  values. These results are compared to the polynomial interpolant of the same data, and the  $\varepsilon \rightarrow 0$  limit is confirmed. This limit cannot be attained with the standard basis, but the stable basis succeeds. The error is averaged over 100 evenly spaced points.

## 2 Existing Parameterization Methods

As we work towards trying to identify a suitable shape parameter, we recall existing techniques in the literature which have had variable levels of success. Each of these existing parameterization techniques has some objective function which needs to be minimized, and that is what we will introduce now. [MJM] *Maybe talk briefly about how a lot of people just kinda guess, or look at small examples to gain intuition, or something? Also, expert knowledge?*

### 2.1 The power function, a.k.a., the Kriging variance

As we saw in (1.5), the error in a scattered data interpolation is bounded by the power function times the Hilbert space norm of  $f$ . Our choice of parameter  $\varepsilon$  will have an effect on  $\|f\|_{\mathcal{H}_K}$ , but it is not one that we can immediately understand. [MJM] *Maybe I'll type more up about this.*

We are, however, able to measure how  $\varepsilon$  affects the power function, and this motivates one strategy that people have used to optimize their kernels: minimize the power function to minimize the pointwise interpolant error. In this setting, the objective function is

$$C_{\text{POWER}}(\varepsilon; k) = \|P_{K, \mathcal{X}}\|_k, \quad (2.1)$$

where the  $k$ -norm can be chosen in one of several ways (2-norm,  $\infty$ -norm,  $\mathcal{H}_K$  norm). If, at the time of interpolation, you know at which points  $\hat{\mathcal{X}}$  you want to evaluate  $s$ , you can choose a discrete norm of  $P_{K, \mathcal{X}}(\hat{x}_i)$ ,  $1 \leq i \leq \hat{N}$ , to minimize the error at those points. If you do not know those points, you can just approximate the  $L^2(\Omega)$  norm.

Historically, the power function has been unstable during computation because of the presence of the  $K^{-1}$  term, which is notoriously ill-conditioned. Using the Hilbert-Schmidt SVD Section 1.1.1 relieves this difficulty because  $\mathbf{k}(\mathbf{x})^T K^{-1} = \boldsymbol{\psi}(\mathbf{x})^T \boldsymbol{\Psi}^{-1}$ :

$$P_{K, \mathcal{X}}(\mathbf{x}) = \sqrt{K(\mathbf{x}, \mathbf{x}) - \boldsymbol{\psi}(\mathbf{x})^T \boldsymbol{\Psi}^{-1} \mathbf{k}(\mathbf{x})}.$$

See Figure 4 for some examples measuring the 2-norm of the power function compared to the error of the interpolation.

In these graphs,  $C_{\text{POWER}}(\varepsilon; 2)$  seems to decrease until reaching roughly  $10^{-8}$ ; this unfortunately seems to be the result of numerical cancelation which occurs as  $\boldsymbol{\psi}(\mathbf{x})^T \boldsymbol{\Psi}^{-1} \mathbf{k}(\mathbf{x}) \rightarrow 1$  for  $\varepsilon \rightarrow 0$ . Fixing this may be possible if we perform additional matrix algebra on  $P_{K, \mathcal{X}}$  using the Hilbert-Schmidt SVD:

$$\begin{aligned} P_{K, \mathcal{X}}(\mathbf{x})^2 &= K(\mathbf{x}, \mathbf{x}) - \mathbf{k}(\mathbf{x})^T K^{-1} \mathbf{k}(\mathbf{x}) \\ &= K(\mathbf{x}, \mathbf{x}) - \boldsymbol{\psi}(\mathbf{x})^T \Lambda_1 \Phi_1^T \Phi_1^{-T} \Lambda_1^{-1} \boldsymbol{\Psi}^{-1} \Phi_1 \Lambda_1 \boldsymbol{\psi}(\mathbf{x}) \\ &= K(\mathbf{x}, \mathbf{x}) - \boldsymbol{\psi}(\mathbf{x})^T \boldsymbol{\Psi}^{-1} \Phi_1 \Lambda_1 \boldsymbol{\psi}(\mathbf{x}) \end{aligned}$$

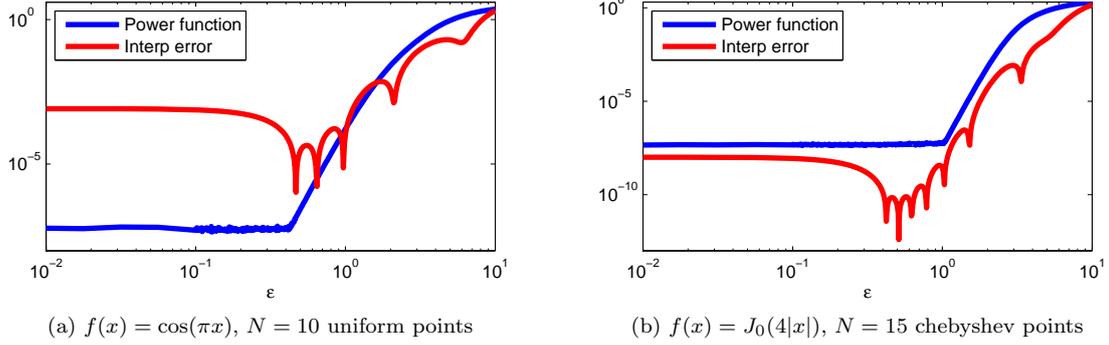


Figure 4: In this example,  $N$  points in  $(-1, 1)$  are sampled and Gaussian interpolation is conducted with a range of  $\varepsilon$  values. These results are compared to  $C_{\text{POWER}}(\varepsilon; 2)$ , which seems to flatten out for small values of  $\varepsilon$ . This suggests that minimizing  $P_{K, \mathcal{X}}$  alone may not be sufficient to obtain an optimal epsilon, though it may help in finding an appropriate region. The error is averaged over 100 evenly spaced points.

To simplify this, we will need to analyze  $\Psi^{-1}\Phi_1\Lambda_1$  recalling (1.9):

$$\begin{aligned}
\Psi^{-1}\Phi_1\Lambda_1 &= \left( (\Phi_1 \quad \Phi_2) \begin{pmatrix} I_N & \\ & \Lambda_2\Phi_2^T\Phi_1^{-T}\Lambda_1^{-1} \end{pmatrix} \right)^{-1} \Phi_1\Lambda_1 \\
&= (\Phi_1 + \Phi_2\Lambda_2\Phi_2^T\Phi_1^{-T}\Lambda_1^{-1})^{-1} \Phi_1\Lambda_1 \\
&= (\Lambda_1^{-1} + \Lambda_1^{-1}\Phi_1^{-1}\Phi_2\Lambda_2\Phi_2^T\Phi_1^{-T}\Lambda_1^{-1})^{-1} \\
&= \Lambda_1^{1/2} \left( I_N + \Lambda_1^{-1/2}\Phi_1^{-1}\Phi_2\Lambda_2\Phi_2^T\Phi_1^{-T}\Lambda_1^{-1/2} \right)^{-1} \Lambda_1^{1/2} \\
&= \Lambda_1^{1/2} (I_N + \mathbf{L}\mathbf{L}^T)^{-1} \Lambda_1^{1/2},
\end{aligned}$$

where we have defined  $\mathbf{L} = \Lambda_2^{1/2}\Phi_2\Phi_1^{-1}\Lambda_1^{-1/2}$ . We want to use the von Neumann series to write

$$(I_N + \mathbf{L}\mathbf{L}^T)^{-1} = \sum_{k=0}^{\infty} (-1)^k (\mathbf{L}\mathbf{L}^T)^k = I_N - \mathbf{L}\mathbf{L}^T + (\mathbf{L}\mathbf{L}^T)^2 - \dots, \quad (2.2)$$

which is only allowed if  $\|\mathbf{L}\mathbf{L}^T\|_2 \leq 1$ . Studying the structure of  $\mathbf{L}$  shows

$$\begin{aligned}
\|\mathbf{L}\mathbf{L}^T\|_2 &\leq \|\mathbf{L}\|_2^2 \leq \|\Lambda_2^{1/2}\|_2^2 \|\Lambda_1^{-1/2}\|_2^2 \|\Phi_2\Phi_1^{-1}\|_2^2 \\
&= (\varepsilon^{N+1})^2 (\varepsilon^{-N})^2 \|\Phi_2\Phi_1^{-1}\|_2^2 = \varepsilon^2 \|\Phi_2\Phi_1^{-1}\|_2^2
\end{aligned}$$

using the submultiplicativity of induced matrix norms. This implies that there is an  $\varepsilon$  small enough such that the series (2.2) can be used. Reinspecting the power function for sufficiently small  $\varepsilon$  with this expansion shows

$$\begin{aligned}
P_{K, \mathcal{X}}(\mathbf{x})^2 &= K(\mathbf{x}, \mathbf{x}) - \boldsymbol{\psi}(\mathbf{x})^T \Lambda_1^{1/2} (I_N - \mathbf{L}\mathbf{L}^T + (\mathbf{L}\mathbf{L}^T)^2 - \dots) \Lambda_1^{1/2} \boldsymbol{\psi}(\mathbf{x}) \\
&= K(\mathbf{x}, \mathbf{x}) - \boldsymbol{\psi}(\mathbf{x})^T \Lambda_1 \boldsymbol{\psi}(\mathbf{x}) + \boldsymbol{\psi}(\mathbf{x})^T \Lambda_1^{1/2} \mathbf{L}\mathbf{L}^T \Lambda_1^{1/2} \boldsymbol{\psi}(\mathbf{x}) - \dots \\
&\approx K(\mathbf{x}, \mathbf{x}) - \boldsymbol{\psi}(\mathbf{x})^T \Lambda_1 \boldsymbol{\psi}(\mathbf{x}),
\end{aligned}$$

where the final approximation becomes more valid as  $\varepsilon \rightarrow 0$ . Because of the swift decay of the eigenvalues it is likely that the terms in  $\boldsymbol{\psi}(\mathbf{x})^T \Lambda_1 \boldsymbol{\psi}(\mathbf{x})$  will need to be separated out to avoid cancelation. Thus we may prefer to write

$$P_{K, \mathcal{X}}(\mathbf{x})^2 \approx K(\mathbf{x}, \mathbf{x}) - \sum_{n=1}^N \lambda_n \psi_n(\mathbf{x})^2 = K(\mathbf{x}, \mathbf{x}) - \lambda_1 \psi_1(\mathbf{x})^2 - \lambda_2 \psi_2(\mathbf{x})^2 - \dots, \quad (2.3)$$

so that the magnitude of the first few terms, which are likely on the same order as  $K(\mathbf{x}, \mathbf{x})$ , do not overwhelm the later terms.

### 2.1.1 The Golomb-Weinberger Error bound

At the end of the Appendix, we suggest that the standard error bound (1.5) could be made computable if our predictive accuracy on  $\Omega$  is good enough:

$$\|f - s\|_{\mathcal{H}_K} \leq \delta_\varepsilon \|s\|_{\mathcal{H}_K}.$$

This term  $\delta_\varepsilon$  is presumably a small number if we have done a good job of approximating  $f$ ; the  $\varepsilon$  term is left there just to show that it is not independent of  $\varepsilon$  as is suggested in Figure 2 for large  $\varepsilon$ .

Using this assumption, and (3.3), we can write (1.5) as

$$|f(\mathbf{x}) - s(\mathbf{x})| \leq \delta_\varepsilon \|s\|_{\mathcal{H}_K} P_{K,\mathcal{X}}(\mathbf{x}), \quad (2.4)$$

which is referred to as the Golomb-Weinberger bound. We determined  $\|s\|_{\mathcal{H}_K} = \sqrt{\mathbf{y}^T \mathbf{K}^{-1} \mathbf{y}}$  in (3.4), so by assuming that we have done a decent job approximating  $f$ , i.e.,  $\|f - s\|_{\mathcal{H}_K}$  is small, our new criterion for optimizing  $\varepsilon$  is

$$C_{\text{GW}}(\varepsilon; k) = \sqrt{\mathbf{y}^T \mathbf{K}^{-1} \mathbf{y}} \|P_{K,\mathcal{X}}\|_k. \quad (2.5)$$

We have already discussed in Section 2.1 how computing the power function can be numerically unstable and requires the Hilbert-Schmidt SVD; we also discussed how  $P_{K,\mathcal{X}}$  is subject to cancelation, but we'll forget that for a moment. The presence of  $\mathbf{K}^{-1}$  in the  $\mathbf{y}^T \mathbf{K}^{-1} \mathbf{y}$  term seems dangerous because of ill-conditioning in  $\mathbf{K}$ . Indeed, direct computation of this term is not recommended, and even with the Hilbert-Schmidt SVD the situation is still troubling. Simple substitution leads to

$$\|s\|_{\mathcal{H}_K}^2 = \mathbf{y}^T \mathbf{K}^{-1} \mathbf{y} = \mathbf{y}^T (\Psi \Lambda_1 \Phi_1^T)^{-1} \mathbf{y} = \mathbf{y}^T \Phi_1^{-T} \Lambda_1^{-1} \Psi^{-1} \mathbf{y} = \mathbf{y}_\Phi^T \Lambda_1^{-1} \mathbf{y}_\Psi, \quad (2.6)$$

where  $\Phi_1 \mathbf{y}_\Phi = \mathbf{y}$  and  $\Psi \mathbf{y}_\Psi = \mathbf{y}$ . Those terms are (presumably) safe to compute, or else the stable basis cannot solve the interpolation problem.

Even so, this term  $\mathbf{y}_\Phi^T \Lambda_1^{-1} \mathbf{y}_\Psi$  may not be safe to compute as it is a weighted inner product between  $\mathbf{y}_\Phi^T$  and  $\mathbf{y}_\Psi$ . This computation is therefore subject to numerical cancelation if the vectors are on different scales; in experiments we have even seen this value computed to be a negative number, despite the fact that  $\|s\|_{\mathcal{H}_K} \geq 0$ .

Although this direct computation has proven unreliable in practice, further manipulations will allow us to bound  $\|s\|_{\mathcal{H}_K}$ , and the bound will become tighter as  $\varepsilon \rightarrow 0$ . The first step is to define  $\Psi \mathbf{b} = \mathbf{y}$  - our goal will be to compute

$$\mathbf{y}^T \mathbf{K}^{-1} \mathbf{y} = \mathbf{b}^T \mathbf{B} \mathbf{b}$$

for some matrix  $\mathbf{B}$ . This vector  $\mathbf{b} = \Psi^{-1} \mathbf{y}$  is safe to compute if  $\Psi$  is indeed a stable basis, and this symmetric inner product will not be subject to the same numerical issues as (2.6). To find  $\mathbf{B}$  we may perform some manipulations involving the HS-SVD  $\mathbf{K} = \Psi \Lambda_1 \Phi_1^T$ ,

$$\begin{aligned} \mathbf{y}^T \mathbf{K}^{-1} \mathbf{y} &= \mathbf{y}^T \Phi_1^{-T} \Lambda_1^{-1} \Psi^{-1} \mathbf{y} \\ &= (\Psi \mathbf{b})^T \Phi_1^{-T} \Lambda_1^{-1} \Psi^{-1} \Psi \mathbf{b} \\ &= \mathbf{b}^T \Psi^T \Phi_1^{-T} \Lambda_1^{-1} \mathbf{b}. \end{aligned} \quad (2.7)$$

Thus  $\mathbf{B} = \Psi^T \Phi_1^{-T} \Lambda_1^{-1}$ , although at first glance this does not seem very helpful.

To simplify  $\Psi^T \Phi_1^{-T} = (\Phi_1^{-1} \Psi)^T$ , we will recall the structure of  $\Psi$  from (1.9):

$$\psi(\mathbf{x})^T = \phi(\mathbf{x})^T \begin{pmatrix} \mathbf{I}_N \\ \Lambda_2 \Phi_2^T \Phi_1^{-T} \Lambda_1^{-1} \end{pmatrix} \Rightarrow \Psi = (\Phi_1 \quad \Phi_2) \begin{pmatrix} \mathbf{I}_N \\ \Lambda_2 \Phi_2^T \Phi_1^{-T} \Lambda_1^{-1} \end{pmatrix},$$

where we have used the fact that  $\Phi = (\Phi_1 \quad \Phi_2)$ . Substituting this into  $\Phi_1^{-1} \Psi$  gives

$$\begin{aligned} \Phi_1^{-1} \Psi &= \Phi_1^{-1} (\Phi_1 \quad \Phi_2) \begin{pmatrix} \mathbf{I}_N \\ \Lambda_2 \Phi_2^T \Phi_1^{-T} \Lambda_1^{-1} \end{pmatrix} \\ &= (\mathbf{I}_N \quad \Phi_1^{-1} \Phi_2) \begin{pmatrix} \mathbf{I}_N \\ \Lambda_2 \Phi_2^T \Phi_1^{-T} \Lambda_1^{-1} \end{pmatrix} \\ &= \mathbf{I}_N + \Phi_1^{-1} \Phi_2 \Lambda_2 \Phi_2^T \Phi_1^{-T} \Lambda_1^{-1}. \end{aligned}$$

Using this result in (2.7) gives

$$\begin{aligned}
\mathbf{y}^T \mathbf{K}^{-1} \mathbf{y} &= \mathbf{b}^T (\Phi_1^{-1} \Psi)^T \Lambda_1^{-1} \mathbf{b} \\
&= \mathbf{b}^T (I_N + \Phi_1^{-1} \Phi_2 \Lambda_2 \Phi_2^T \Phi_1^{-T} \Lambda_1^{-1})^T \Lambda_1^{-1} \mathbf{b} \\
&= \mathbf{b}^T \Lambda_1^{-1} \mathbf{b} + \mathbf{b}^T \Lambda_1^{-1} \Phi_1^{-1} \Phi_2 \Lambda_2 \Phi_2^T \Phi_1^{-T} \Lambda_1^{-1} \mathbf{b}
\end{aligned} \tag{2.8}$$

$$\begin{aligned}
&= \mathbf{b}^T \Lambda_1^{-1} \mathbf{b} + (\Lambda_2^{1/2} \Phi_2^T \Phi_1^{-T} \Lambda_1^{-1} \mathbf{b})^T (\Lambda_2^{1/2} \Phi_2^T \Phi_1^{-T} \Lambda_1^{-1} \mathbf{b}) \\
&\geq \mathbf{b}^T \Lambda_1^{-1} \mathbf{b},
\end{aligned} \tag{2.9}$$

where the inequality is true because  $\|\Lambda_2^{1/2} \Phi_2^T \Phi_1^{-T} \Lambda_1^{-1} \mathbf{b}\|_2^2 \geq 0$ . Here,  $\Lambda_2^{1/2}$  is defined as the diagonal matrix with the positive square roots of  $\Lambda_2$  on its diagonal.

This result implies that if we use the Hilbert-Schmidt SVD we can bound the native space norm of the interpolant from below by  $\sqrt{\mathbf{b}^T \Lambda_1^{-1} \mathbf{b}}$ . This term is safely-computable because it is a symmetric inner product, meaning that cancelation will not occur:

$$\mathbf{b}^T \Lambda_1^{-1} \mathbf{b} = \sum_{n=1}^N \frac{b_n^2}{\lambda_n}.$$

We can reach this result another way, by instead studying  $\|s\|_{\mathcal{H}_K}^2$  with the Hilbert-Schmidt SVD (recalling (1.11) and (1.9))

$$\begin{aligned}
\|s\|_{\mathcal{H}_K}^2 &= \langle s^T, s \rangle_{\mathcal{H}_K} \\
&= \langle \mathbf{y}^T \Psi^T \psi(\cdot), \psi(\cdot)^T \Psi \mathbf{y} \rangle_{\mathcal{H}_K} \\
&= \mathbf{y}^T \Psi^T \left\langle \left( \Lambda_2 \Phi_2^T \Phi_1^{-T} \Lambda_1^{-1} \right)^T \phi(\cdot), \phi(\cdot)^T \left( \Lambda_2 \Phi_2^T \Phi_1^{-T} \Lambda_1^{-1} \right) \right\rangle_{\mathcal{H}_K} \Psi \mathbf{y} \\
&= \mathbf{b}^T (I_N \quad \Lambda_1^{-1} \Phi_1^{-1} \Phi_2 \Lambda_2) \langle \phi(\cdot), \phi(\cdot)^T \rangle_{\mathcal{H}_K} \begin{pmatrix} I_N \\ \Lambda_2 \Phi_2^T \Phi_1^{-T} \Lambda_1^{-1} \end{pmatrix} \mathbf{b} \\
&= \mathbf{b}^T (I_N \quad \Lambda_1^{-1} \Phi_1^{-1} \Phi_2 \Lambda_2) \begin{pmatrix} \Lambda_1^{-1} & \\ & \Lambda_2^{-1} \end{pmatrix} \begin{pmatrix} I_N \\ \Lambda_2 \Phi_2^T \Phi_1^{-T} \Lambda_1^{-1} \end{pmatrix} \mathbf{b} \\
&= \mathbf{b}^T (\Lambda_1^{-1} + \Lambda_1^{-1} \Phi_1^{-1} \Phi_2 \Lambda_2 \Phi_2^T \Phi_1^{-T} \Lambda_1^{-1}) \mathbf{b}.
\end{aligned}$$

The resolution of the inner product is a result of the Hilbert space orthogonality of the eigenfunctions described in (1.8), i.e.,  $\langle \phi(\cdot), \phi(\cdot)^T \rangle_{\mathcal{H}_K} = \Lambda^{-1}$ .

**Remark 1** *Much of what is presented in this paper applied to general kernels, but the lower bound on  $\|s\|_{\mathcal{H}_K}^2$  derived below applies only to Gaussian kernels. There are likely similar bounds on other kernels as well. A thorough discussion of the Gaussian eigenvalues is available in [3].*

The bound (2.9) becomes tighter as  $\varepsilon \rightarrow 0$ , because the remainder term in (2.8) is orders of magnitude smaller than the main term. [GEF] [There are pictures to illustrate this in my Sicily talk.](#) To prove this, we need to study  $\|\Lambda_2^{1/2} \Phi_2^T \Phi_1^{-T} \Lambda_1^{-1} \mathbf{b}\|_2^2$ ; we assume that  $\Phi_2^T \Phi_1^{-T}$  is a well-behaved matrix as  $\varepsilon \rightarrow 0$ . This is a reasonable assumption if the eigenfunctions are properly scaled, a topic which is discussed in [3]. Furthermore, we can note that  $\mathbf{b}$  is well-behaved as  $\varepsilon \rightarrow 0$  because otherwise  $\psi$  would not be a stable basis.

We need to consider the magnitude of terms in  $\Lambda_2^{1/2} \Phi_2^T \Phi_1^{-T} \Lambda_1^{-1}$  as  $\varepsilon \rightarrow 0$ , under the assumption that  $\Phi_2^T \Phi_1^{-T}$  is at its limit. Rows of  $\Phi_2^T \Phi_1^{-T}$  are scaled by  $\Lambda_2^{1/2}$  and columns are scaled by  $\Lambda_1^{-1}$ . Note that as  $\varepsilon \rightarrow 0$ ,

$$(\Lambda_1^{-1})_{ii} \approx \varepsilon^{-2i}, \quad (\Lambda_2^{1/2})_{ii} \approx \varepsilon^{i+N}. \tag{2.10}$$

Recalling a property of matrix norms that I found on Wikipedia (I'm sure it's also in [5]), we know that  $\|\mathbf{A}\mathbf{x}\|_a \leq \|\mathbf{A}\|_b \|\mathbf{x}\|_a$  for  $a, b$  such that  $\|\cdot\|_b$  is an induced matrix norm and  $\|\cdot\|_a$  is a valid vector norm. Using that, (2.10), and

the submultiplicativity of induced matrix norms, we know

$$\begin{aligned} \lim_{\varepsilon \rightarrow 0} \left\| \Lambda_2^{1/2} \Phi_2^T \Phi_1^{-T} \Lambda_1^{-1} \mathbf{b} \right\|_2 &\leq \lim_{\varepsilon \rightarrow 0} \left\| \Lambda_2^{1/2} \Phi_2^T \Phi_1^{-T} \Lambda_1^{-1} \right\|_\infty \|\mathbf{b}\|_2 \\ &\leq \|\mathbf{b}\|_2 \|\Phi_2^T \Phi_1^{-T}\|_\infty \lim_{\varepsilon \rightarrow 0} \left\| \Lambda_2^{1/2} \right\|_\infty \|\Lambda_1^{-1}\|_\infty \\ &\leq \|\mathbf{b}\|_2 \|\Phi_2^T \Phi_1^{-T}\|_\infty \varepsilon^{1+N} \varepsilon^{-2N} \end{aligned}$$

This lets us say that

$$\lim_{\varepsilon \rightarrow 0} \left\| \Lambda_2^{1/2} \Phi_2^T \Phi_1^{-T} \Lambda_1^{-1} \mathbf{b} \right\|_2^2 \leq \varepsilon^{2-2N} \|\Phi_2^T \Phi_1^{-T}\|_\infty^2 \|\mathbf{b}\|_2^2.$$

When we compare that to

$$\mathbf{b}^T \Lambda_1^{-1} \mathbf{b} = \left\| \Lambda_1^{-1/2} \mathbf{b} \right\|_2^2 \leq \left\| \Lambda_1^{-1/2} \right\|_\infty^2 \|\mathbf{b}\|_2^2 = (\varepsilon^{-N})^2 \|\mathbf{b}\|_2^2 = \varepsilon^{-2N} \|\mathbf{b}\|_2^2,$$

we see that, as  $\varepsilon \rightarrow 0$ , the remainder term of (2.8) grows two orders of magnitude more slowly than the bound in (2.9).

[GEF] So if you combine the two ideas, does  $C_{GW}(\varepsilon; k)$  give you a good criterion? I've used the  $C_{GW}(\varepsilon; \infty)$  criterion—without the HS-SVD—before. See Fig. 20 in the past, present, future paper. For me, it seemed to be right there with MLE. [MJM] I think we're still subject to the same cancelation for the power function. The fix I just typed up may correct that. I'll need to test it though.

## 2.2 Cross-Validation

This writeup is stolen liberally from the Hickernell writeup that I mentioned earlier. The idea of cross-validation is essentially:

- Given the scattered data  $\mathcal{X}$  and  $\mathbf{y}$ , partition the design into disjoint, nonempty groups  $\mathcal{X}_I$  and  $\mathcal{X}_O$  so that  $\mathcal{X}_I \cup \mathcal{X}_O = \mathcal{X}$ . The set  $\mathcal{X}_I$  will be used to create an auxiliary approximation  $s_I$  and the set  $\mathcal{X}_O$  will be used to judge the accuracy of  $s_I$ .

– The sets  $\mathcal{X}_I$  and  $\mathcal{X}_O$  block up the interpolation matrix  $\mathbf{K}$ , its inverse  $\mathbf{A}$ , and the vectors  $\mathbf{y}$  and  $\mathbf{c}$  as

$$\mathbf{K} = \begin{pmatrix} \mathbf{K}_{II} & \mathbf{K}_{IO} \\ \mathbf{K}_{OI} & \mathbf{K}_{OO} \end{pmatrix}, \quad \mathbf{A} = \begin{pmatrix} \mathbf{A}_{II} & \mathbf{A}_{IO} \\ \mathbf{A}_{OI} & \mathbf{A}_{OO} \end{pmatrix}, \quad \mathbf{y} = \begin{pmatrix} \mathbf{y}_I \\ \mathbf{y}_O \end{pmatrix}, \quad \mathbf{c} = \begin{pmatrix} \mathbf{c}_I \\ \mathbf{c}_O \end{pmatrix}.$$

Recall that  $\mathbf{K}\mathbf{c} = \mathbf{y}$ , and thus  $\mathbf{c} = \mathbf{A}\mathbf{y}$ ; because  $\mathbf{K}$  and  $\mathbf{A}$  are symmetric, we know that  $\mathbf{K}_{IO} = \mathbf{K}_{OI}^T$  and  $\mathbf{A}_{IO} = \mathbf{A}_{OI}^T$ . Also recall that the parameter  $\varepsilon$  appears in  $\mathbf{K}$ .

- The available data at  $\mathcal{X}_0$  is  $\mathbf{y}_0$ , and the prediction at the points  $\mathcal{X}_0$  using the data  $(\mathcal{X}_I, \mathbf{y}_I)$  is  $\mathbf{K}_{OI} \mathbf{K}_{II}^{-1} \mathbf{y}_I$  which can be determined by applying the structure of (1.4) at  $x$  locations in  $\mathcal{X}_O$ . Therefore, studying

$$|\mathbf{y}_O - \mathbf{K}_{OI} \mathbf{K}_{II}^{-1} \mathbf{y}_I|$$

tells us something about how good our approximation is. The matrix  $\mathbf{K}_{OI} \mathbf{K}_{II}^{-1}$  can be thought of as an operator which takes in values on  $\mathcal{X}_I$  and interpolates them to  $\mathcal{X}_O$  [6].

- By considering a set of  $p$  partitions  $\mathcal{O} = \{\mathcal{X}_O^{(1)}, \dots, \mathcal{X}_O^{(p)}\}$  such that

$$\mathcal{X}_O^{(i)} \cap \mathcal{X}_O^{(j)} = \emptyset, i \neq j, \quad \text{and} \quad \bigcup_{i=1}^p \mathcal{X}_O^{(i)} = \mathcal{X}$$

with  $\mathcal{X}_I^{(i)} = \mathcal{X} \setminus \mathcal{X}_O^{(i)}$ , we could consider instances with certain pieces of data omitted for certain  $i$  values. Then we could consider the residual left-over by our interpolants evaluated at those points:

$$C_{CV}(\varepsilon; \mathcal{O}, k) = \sum_{\mathcal{X}_O \in \mathcal{O}} \|\mathbf{y}_O - \mathbf{K}_{OI} \mathbf{K}_{II}^{-1} \mathbf{y}_I\|_k, \quad (2.11)$$

where  $k$  is probably 1, 2 or  $\infty$ .

The criterion to be minimized for a good  $\varepsilon$  value is (2.11). Of course, the presence of  $\mathbf{K}_{II}^{-1}$  suggests ill-conditioning may be a problem. If we choose to implement the Hilbert-Schmidt SVD here, we should block up the interpolation matrix

$$\mathbf{K} = \begin{pmatrix} \mathbf{K}_{II} & \mathbf{K}_{IO} \\ \mathbf{K}_{OI} & \mathbf{K}_{OO} \end{pmatrix} = \begin{pmatrix} \Psi_{II}\Lambda_I\Phi_I^T & \Psi_{IO}\Lambda_O\Phi_O^T \\ \Psi_{OI}\Lambda_I\Phi_I^T & \Psi_{OO}\Lambda_O\Phi_O^T \end{pmatrix} = \begin{pmatrix} \Psi_{II} & \Psi_{IO} \\ \Psi_{OI} & \Psi_{OO} \end{pmatrix} \begin{pmatrix} \Lambda_I\Phi_I^T & \\ & \Lambda_O\Phi_O^T \end{pmatrix}. \quad (2.12)$$

Using this block structure, we see that  $\mathbf{K}_{II} = \Psi_{II}\Lambda_I\Phi_I^T$  and  $\mathbf{K}_{OI} = \Psi_{OI}\Lambda_I\Phi_I^T$ , so therefore we can write the criterion as

$$C_{CV}(\varepsilon; \mathcal{O}, k) = \sum_{\mathbf{x}_o \in \mathcal{O}} \|\mathbf{y}_o - \Psi_{OI}\Psi_{II}^{-1}\mathbf{y}_I\|_k, \quad (2.13)$$

which is presumably safer to compute.

The expression in (2.11) can be rewritten in a slightly different way by exploiting the block structure of the relevant vectors and matrices. Note that  $\mathbf{c} = \mathbf{A}\mathbf{y}$  implies

$$\mathbf{c}_o = \mathbf{A}_{OI}\mathbf{y}_I + \mathbf{A}_{OO}\mathbf{y}_o.$$

We can also invoke the fact that  $\mathbf{A}\mathbf{K} = \mathbf{I}_N$  (where  $\mathbf{I}_N$  is the  $N \times N$  identity) to note that

$$\mathbf{A}_{OI}\mathbf{K}_{II} + \mathbf{A}_{OO}\mathbf{K}_{IO}^T = 0 \quad \Rightarrow \quad \mathbf{A}_{OI} = -\mathbf{A}_{OO}\mathbf{K}_{IO}^T\mathbf{K}_{II}^{-1}.$$

Plugging this in above gives

$$\begin{aligned} \mathbf{c}_o &= -\mathbf{A}_{OO}\mathbf{K}_{OI}\mathbf{K}_{II}^{-1}\mathbf{y}_I + \mathbf{A}_{OO}\mathbf{y}_o \\ \mathbf{A}_{OO}^{-1}\mathbf{c}_o &= -\mathbf{K}_{OI}\mathbf{K}_{II}^{-1}\mathbf{y}_I + \mathbf{y}_o \end{aligned}$$

which allows us to write our optimization criterion as

$$C_{CV}(\varepsilon; \mathcal{O}, k) = \sum_{\mathbf{x}_o \in \mathcal{O}} \|\mathbf{A}_{OO}^{-1}\mathbf{c}_o\|_k. \quad (2.14)$$

To involve the Hilbert-Schmidt SVD, we need to study (2.12) and define

$$\begin{pmatrix} \Psi_{II} & \Psi_{IO} \\ \Psi_{OI} & \Psi_{OO} \end{pmatrix} \begin{pmatrix} \mathbf{B}_{II} & \mathbf{B}_{IO} \\ \mathbf{B}_{OI} & \mathbf{B}_{OO} \end{pmatrix} = \begin{pmatrix} \mathbf{I}_{N_I} & \\ & \mathbf{I}_{N_O} \end{pmatrix},$$

where  $\mathcal{X}_I$  and  $\mathcal{X}_O$  have  $N_I$  and  $N_O$  points in them respectively. Manipulations similar to those above show that

$$\mathbf{B}_{OO}^{-1} = \Psi_{OO} - \Psi_{OI}\Psi_{II}^{-1}\Psi_{IO},$$

which we will use shortly. We also need to define

$$\begin{pmatrix} \Psi_{II} & \Psi_{IO} \\ \Psi_{OI} & \Psi_{OO} \end{pmatrix} \begin{pmatrix} \mathbf{b}_I \\ \mathbf{b}_O \end{pmatrix} = \begin{pmatrix} \mathbf{y}_I \\ \mathbf{y}_O \end{pmatrix} \quad \Rightarrow \quad \begin{pmatrix} \mathbf{b}_I \\ \mathbf{b}_O \end{pmatrix} = \begin{pmatrix} \mathbf{B}_{II} & \mathbf{B}_{IO} \\ \mathbf{B}_{OI} & \mathbf{B}_{OO} \end{pmatrix} \begin{pmatrix} \mathbf{y}_I \\ \mathbf{y}_O \end{pmatrix}, \quad \begin{pmatrix} \mathbf{b}_I \\ \mathbf{b}_O \end{pmatrix} = \begin{pmatrix} \Lambda_I\Phi_I^T & \\ & \Lambda_O\Phi_O^T \end{pmatrix} \begin{pmatrix} \mathbf{c}_I \\ \mathbf{c}_O \end{pmatrix}$$

This implies that, in terms of the HS-SVD system,  $\mathbf{c}_O = \Phi_O^{-T}\Lambda_O^{-1}\mathbf{b}_O$ , and sufficient study of (2.12) knowing that  $\mathbf{A}\mathbf{K} = \mathbf{I}_N$  gives us  $\mathbf{A}_{OO} = \Phi_O^{-T}\Lambda_O^{-1}\mathbf{B}_{OO}$ . Combining these in (2.14) gives us

$$C_{CV}(\varepsilon; \mathcal{O}, k) = \sum_{\mathbf{x}_o \in \mathcal{O}} \|\mathbf{B}_{OO}^{-1}\mathbf{b}_O\|_k = \sum_{\mathbf{x}_o \in \mathcal{O}} \|(\Psi_{OO} - \Psi_{OI}\Psi_{II}^{-1}\Psi_{IO})^{-1}\mathbf{b}_O\|_k, \quad (2.15)$$

which is almost certainly more stable to compute.

Often times, cross-validation is conducted in one of two ways:

- **Leave-one-out cross-validation** - All the data except a single point is used to compute the interpolant, and the residual is judged at that point. In this setting,  $\mathcal{O} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$  and the errors at each of those points are added up to find  $CV(\varepsilon; \mathcal{O}, k)$ . (2.14) may be preferable to compute in this case because  $\mathbf{A}_{OO}$  is just a number.

- **Leave-half-out cross-validation** - Half of the data is omitted to create an interpolant and the residual is judged on the other half; then the process is flipped and both results are combined to compute  $CV(\varepsilon; \mathcal{O}, k)$ . In this setting,  $\mathcal{O} = \{\mathcal{X}_O^{(1)}, \mathcal{X}_O^{(2)}\}$  and  $|\mathcal{X}_O^{(1)}| = |\mathcal{X}_O^{(2)}|$ , or as close as possible. (2.11) is almost surely the preferred computation in this case because  $A_{OO}$  is roughly size  $N/2 \times N/2$  and computing its inverse is costly.

An example of cross-validation results is presented in Figure 5, although the results are somewhat inconclusive: the minimum in Figure 5a seems near the “true”  $\varepsilon$ , but in Figure 5b the  $C_{CV}$  function for LOOCV seems to drop for  $\varepsilon \rightarrow 0$  which does not reflect the error curve.

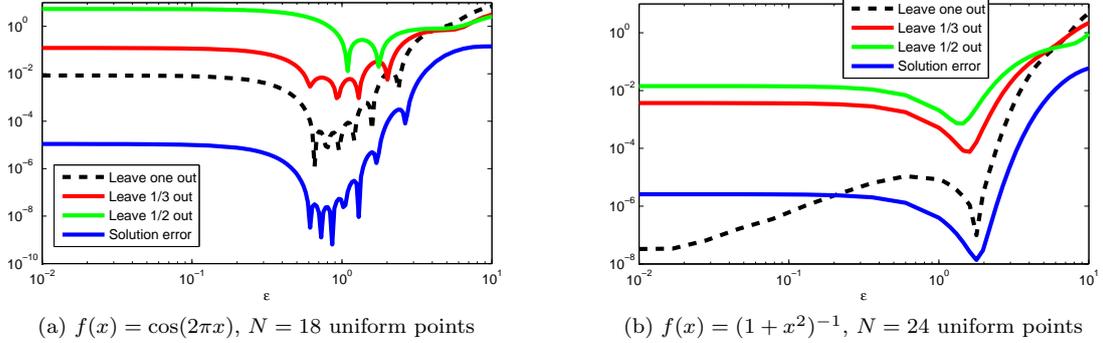


Figure 5: In this example,  $N$  points in  $(-1, 1)$  are sampled and Gaussian interpolation is conducted with a range of  $\varepsilon$  values. These results are compared to  $C_{CV}(\varepsilon; \mathcal{O}, 2)$ , for three  $\mathcal{O}$  choices: one with 2 equally sized sets, 3 equally sized set, and  $N$  equally sized sets. The error is averaged over 100 evenly spaced points.

One concern, which is discussed in Section 2.3.1, is the process variance. This is the idea that the kernel  $K$  in the Gaussian field could actually be written as  $\hat{K} = \sigma K$  for some positive valued  $\sigma$  called the *process variance*. This plays no role here because such a term gets canceled out in the computation of  $K_{IO}^T K_{II}^{-1} \mathbf{y}_I$ , in the same way that it is canceled out in  $s(\mathbf{x}) = \mathbf{k}(\mathbf{x})^T K^{-1} \mathbf{y}$ . It does potentially play a role in other settings though, most notably the Kriging variance, and is discussed later.

Although often described in a statistical setting, this cross-validation concept is equally well supported from the numerical analysis side. No rigorous structure that I know of exists to explain why this is a useful technique analytically, as opposed to just an idea that makes sense.

### 2.3 Maximum Likelihood Estimation

This approach leans primarily on the Gaussian processes framework, though we will show later how it can be derived using the Hilbert space framework. Begin by recalling (1.13), that  $\mathbf{Y}_{\mathcal{X}} \sim \mathcal{N}(\boldsymbol{\mu}, \mathbf{K})$ , and, as before, restricting the Gaussian Process  $Y$  to have zero mean, i.e.,  $\boldsymbol{\mu}_Y \equiv 0$ . We are going to treat  $\varepsilon$  as *drawn from a random variable*  $\mathcal{E}$  now, with some unknown distribution; in turn we also will need to define the joint random variable  $Z = (\mathcal{E}, \mathbf{Y}_{\mathcal{X}})$  with density  $p_Z$ . We want to study the conditional density function  $p_{\mathcal{E}|\mathbf{Y}_{\mathcal{X}}}(\varepsilon|\mathbf{Y}_{\mathcal{X}} = \mathbf{y})$ , that is, how likely it is that the Gaussian process had covariance  $K$ , parameterized by  $\varepsilon$ , when realizing the data  $\mathbf{y}$ .

This function is often called the *likelihood* function, and it allows us to compare the relative likelihood of values of  $\varepsilon$  given the existing data. Maximizing this function yields (in some sense) the  $\varepsilon$  which most likely parameterized the covariance kernel which generated the data  $\mathbf{y}$ . The function  $p_{\mathcal{E}|\mathbf{Y}_{\mathcal{X}}}(\varepsilon|\mathbf{Y}_{\mathcal{X}} = \mathbf{y})$  is a function of  $\varepsilon$  for any fixed value of  $\mathbf{y}$ .

A technical note: the notation from Section 1.2 and this section will not perfectly align. In that section we considered  $\varepsilon$  fixed and the joint distribution  $(\mathbf{Y}_{\mathcal{X}}, Y_{\mathbf{x}_0})$ , but here we are studying the joint distribution  $(\mathcal{E}, \mathbf{Y}_{\mathcal{X}})$ . Where we previously had written  $p_{\mathbf{Y}_{\mathcal{X}}}(\mathbf{y})$  in (1.16), the appropriate notation in this section is

$$p_{\mathbf{Y}_{\mathcal{X}}|\mathcal{E}}(\mathbf{y}|\mathcal{E} = \varepsilon) = \frac{1}{\sqrt{2\pi \det(\mathbf{K})}} \exp\left(-\frac{1}{2}\mathbf{y}^T \mathbf{K}^{-1} \mathbf{y}\right), \quad (2.16)$$

where although  $\varepsilon$  does not explicitly appear on the right hand side, it appears within  $\mathbf{K}$ .

We want to study  $p_{\varepsilon|\mathbf{Y}_{\mathcal{X}}}(\varepsilon|\mathbf{Y}_{\mathcal{X}} = \mathbf{y})$ , but we do not know that density. What we do know is the definition of conditional probability:

$$P(A|B) = \frac{P(A \cap B)}{P(B)} \Leftrightarrow P(A \cap B) = P(A|B)P(B),$$

assuming  $P(B) > 0$ . Using an analogous version of this for density functions instead of the probability function allows us to write

$$\begin{aligned} p_{\varepsilon|\mathbf{Y}_{\mathcal{X}}}(\varepsilon|\mathbf{Y}_{\mathcal{X}} = \mathbf{y}) &= \frac{p_Z(\varepsilon, \mathbf{y})}{p_{\mathbf{Y}_{\mathcal{X}}}(\mathbf{Y}_{\mathcal{X}} = \mathbf{y})} \\ &= \frac{p_{\mathbf{Y}_{\mathcal{X}}|\varepsilon}(\mathbf{y}|\varepsilon)p_{\varepsilon}(\varepsilon)}{p_{\mathbf{Y}_{\mathcal{X}}}(\mathbf{y})}. \end{aligned} \quad (2.17)$$

We can clean up our likelihood function expression somewhat by ignoring things we have no knowledge of. First off, we do not know what the marginal distribution of  $\varepsilon$  is - if we did we would just study that to determine an optimal  $\varepsilon$  parameterization. We also do not know what the marginal distribution of  $\mathbf{Y}_{\mathcal{X}}$  is: we do not know how  $\mathbf{Y}_{\mathcal{X}}$  varies independently of  $\varepsilon$ . To determine this, we would need to compute

$$p_{\mathbf{Y}_{\mathcal{X}}}(\mathbf{y}) = \int_0^{\infty} p_Z(\mathbf{y}, \varepsilon) d\varepsilon$$

which is not going to happen because the joint density  $p_Z$  is unknown. More importantly,  $p_{\mathbf{Y}_{\mathcal{X}}}(\mathbf{y})$  is independent of  $\varepsilon$ , and therefore changing  $\varepsilon$  will have no effect on that value.

By abandoning  $p_{\varepsilon}(\varepsilon)$  and  $p_{\mathbf{Y}_{\mathcal{X}}}(\mathbf{y})$ , we can suggest that

$$p_{\varepsilon|\mathbf{Y}_{\mathcal{X}}}(\varepsilon|\mathbf{Y}_{\mathcal{X}} = \mathbf{y}) \propto p_{\mathbf{Y}_{\mathcal{X}}|\varepsilon}(\mathbf{y}|\varepsilon), \quad (2.18)$$

and we know this function (from (2.16)). Thus the concept of maximizing the likelihood requires maximizing  $p_{\mathbf{Y}_{\mathcal{X}}|\varepsilon}(\mathbf{y}|\varepsilon)$ . By itself, this function is subject to overflow and underflow, so it is common to instead work with its logarithm:

$$\log(p_{\mathbf{Y}_{\mathcal{X}}|\varepsilon}(\mathbf{y}|\varepsilon)) = -\frac{1}{2} \log \det \mathbf{K} - \frac{1}{2} \mathbf{y}^T \mathbf{K}^{-1} \mathbf{y} - \frac{1}{2} \log 2\pi.$$

Because we have been in the practice of minimizing functions to find optimal  $\varepsilon$  parameterizations, we will multiply by  $-2$  and ignore the constant  $\log 2\pi$  to create our maximum likelihood criterion

$$\begin{aligned} C_{\text{MLE}}(\varepsilon) &= -2 \log(p_{\mathbf{Y}_{\mathcal{X}}|\varepsilon}(\mathbf{y}|\varepsilon)) - \log 2\pi \\ &= \log \det \mathbf{K} + \mathbf{y}^T \mathbf{K}^{-1} \mathbf{y} \end{aligned} \quad (2.19)$$

### 2.3.1 Incorporating the process variance

Mentioned briefly in Section 2.2 was the concept of a process variance  $\sigma > 0$  which defines the maximum covariance of  $K$ ; in this setting, our standard kernel  $K$  would actually be multiplied by  $\sigma$  to create the kernel  $\tilde{K} = \sigma K$  which defines the Gaussian process  $Y \sim GF(\mu_Y, \tilde{K})$ .

Thus far we have considered only  $\sigma = 1$ , and in fact this is acceptable for prediction purposes: recall that  $s(\mathbf{x}) = \mathbf{k}(\mathbf{x})^T \mathbf{K}^{-1} \mathbf{y}$ . Replacing  $K$  with  $\tilde{K}$  would give

$$s(\mathbf{x}) = \tilde{\mathbf{k}}(\mathbf{x})^T \tilde{\mathbf{K}}^{-1} \mathbf{y} = \sigma \mathbf{k}(\mathbf{x})^T (\sigma \mathbf{K})^{-1} \mathbf{y} = \mathbf{k}(\mathbf{x})^T \mathbf{K}^{-1} \mathbf{y}$$

thus the interpolant would be the same for any  $\sigma$ . The same computation can be applied to  $C_{\text{CV}}$  to show that it is invariant for different  $\sigma$  values. The power function  $P_{\tilde{K}, \mathcal{X}} \rightarrow 0$  as  $\sigma \rightarrow 0$ ,

$$P_{\tilde{K}, \mathcal{X}}(\mathbf{x}) = \sqrt{\tilde{K}(\mathbf{x}, \mathbf{x}) - \tilde{\mathbf{k}}(\mathbf{x})^T \tilde{\mathbf{K}}^{-1} \tilde{\mathbf{k}}(\mathbf{x})} = \sqrt{\sigma K(\mathbf{x}, \mathbf{x}) - \sigma \mathbf{k}(\mathbf{x})^T (\sigma \mathbf{K})^{-1} \sigma \mathbf{k}(\mathbf{x})} = \sqrt{\sigma} P_{K, \mathcal{X}}(\mathbf{x}),$$

so when studying the power function it is necessary to fix  $\sigma$ , probably to 1, and vary  $\varepsilon$  to minimize  $C_{\text{POWER}}$ . If the term  $\|f\|_{\mathcal{H}_K}$  or  $\|s\|_{\mathcal{H}_K}$  is included to compute the Native space norm error or Golomb-Weinberger bound, respectively, (see the Appendix) then the situation becomes more complicated because  $\sigma$  is involved in the Native space norm.

The situation is also tricky for maximum likelihood estimation, which is why it is under consideration right now. When we introduce a process variance, we are suggesting that  $\sigma$  is a draw from a random variable  $\Sigma$  with unknown distribution. We would need to study the joint distribution  $(\Sigma, \mathcal{E}, \mathbf{Y}_{\mathcal{X}})$ , and our kernel parameterization would require optimizing for both  $\sigma$  and  $\varepsilon$  by maximizing  $p_{\Sigma, \mathcal{E} | \mathbf{Y}_{\mathcal{X}}}(\sigma, \varepsilon | \mathbf{Y}_{\mathcal{X}} = \mathbf{y})$ , which we will maximize by maximizing  $p_{\mathbf{Y}_{\mathcal{X}} | \Sigma, \mathcal{E}}(\mathbf{y} | \Sigma = \sigma, \mathcal{E} = \varepsilon)$  similarly to (2.18).

We could treat this as a two dimensional optimization problem, but instead we will invoke the technique of *profile likelihood*, where  $\sigma$  will be defined as a function of  $\varepsilon$ , i.e.,  $\sigma \equiv \sigma(\varepsilon)$ . Our goal now is to choose an optimal process variance  $\sigma_{\text{opt}}$  which we will do by maximizing  $p_{\Sigma | \mathcal{E}, \mathbf{Y}_{\mathcal{X}}}(\sigma | \mathcal{E} = \varepsilon, \mathbf{Y}_{\mathcal{X}} = \mathbf{y})$ . The term profile likelihood is a bit of a misnomer, because  $p_{\Sigma, \mathcal{E} | \mathbf{Y}_{\mathcal{X}}}(\sigma(\varepsilon), \varepsilon | \mathbf{Y}_{\mathcal{X}} = \mathbf{y})$  is not derived from a cumulative distribution function and thus it is not a true density and loses some desirable properties. Even so, this is a common technique.

Using the same proportionality logic as in (2.18) we can write that

$$\begin{aligned} p_{\Sigma | \mathcal{E}, \mathbf{Y}_{\mathcal{X}}}(\sigma | \mathcal{E} = \varepsilon, \mathbf{Y}_{\mathcal{X}} = \mathbf{y}) &\propto p_{\mathcal{E}, \mathbf{Y}_{\mathcal{X}} | \Sigma}(\varepsilon, \mathbf{y} | \Sigma = \sigma) \\ &= p_{\mathbf{Y}_{\mathcal{X}} | \Sigma, \mathcal{E}}(\mathbf{y} | \Sigma = \sigma, \mathcal{E} = \varepsilon) p_{\mathcal{E} | \Sigma}(\varepsilon | \Sigma = \sigma) \\ &\propto p_{\mathbf{Y}_{\mathcal{X}} | \Sigma, \mathcal{E}}(\mathbf{y} | \Sigma = \sigma, \mathcal{E} = \varepsilon). \end{aligned}$$

Therefore, our optimal  $\sigma$  can be found by maximizing  $p_{\mathbf{Y}_{\mathcal{X}} | \Sigma, \mathcal{E}}(\mathbf{y} | \Sigma = \sigma, \mathcal{E} = \varepsilon)$ ; this function is the same as (2.16), except with  $\mathbf{K}$  replaced by  $\tilde{\mathbf{K}}$ . As before, instead of maximizing, we will try to minimize the negative log of this function:

$$\begin{aligned} -2 \log (p_{\mathbf{Y}_{\mathcal{X}} | \Sigma, \mathcal{E}}(\mathbf{y} | \Sigma = \sigma, \mathcal{E} = \varepsilon)) + \log 2\pi &= \log \det \tilde{\mathbf{K}} + \mathbf{y}^T \tilde{\mathbf{K}}^{-1} \mathbf{y} \\ &= \log \det \sigma \mathbf{K} + \mathbf{y}^T (\sigma \mathbf{K})^{-1} \mathbf{y} \\ &= N \log \sigma + \log \det \mathbf{K} + \frac{1}{\sigma} \mathbf{y}^T \mathbf{K}^{-1} \mathbf{y} \end{aligned}$$

Differentiating this with respect to  $\sigma$ , setting it equal to 0, and solving for  $\sigma$  gives the optimal profile variance

$$\sigma_{\text{opt}} = \frac{1}{N} \mathbf{y}^T \mathbf{K} \mathbf{y}. \quad (2.20)$$

Using the profile likelihood strategy, we maximize  $p_{\Sigma, \mathcal{E} | \mathbf{Y}_{\mathcal{X}}}(\sigma, \varepsilon | \mathbf{Y}_{\mathcal{X}} = \mathbf{y})$  by minimizing

$$\begin{aligned} -2 \log (p_{\mathbf{Y}_{\mathcal{X}} | \Sigma, \mathcal{E}}(\mathbf{y} | \Sigma = \sigma_{\text{opt}}, \mathcal{E} = \varepsilon)) + \log 2\pi &= N \log \left( \frac{1}{N} \mathbf{y}^T \mathbf{K}^{-1} \mathbf{y} \right) + \log \det \mathbf{K} + \left( \frac{1}{N} \mathbf{y}^T \mathbf{K}^{-1} \mathbf{y} \right)^{-1} \mathbf{y}^T \mathbf{K}^{-1} \mathbf{y} \\ &= N \log (\mathbf{y}^T \mathbf{K}^{-1} \mathbf{y}) + \log \det \mathbf{K} - N \log N + N \end{aligned}$$

thus defining our profile likelihood  $\varepsilon$  parameterization criterion as

$$C_{\text{MPLE}}(\varepsilon) = N \log (\mathbf{y}^T \mathbf{K}^{-1} \mathbf{y}) + \log \det \mathbf{K} \quad (2.21)$$

after omitting the constant  $-N \log N + N$ .

### 2.3.2 A deterministic derivation of MLE

The title of this section is a bit misleading, because of course likelihoods cannot be discussed outside of a probabilistic setting. The criterion  $C_{\text{MPLE}}$  is equivalent to a criterion which can be derived deterministically, which we will show now. This too was swiped from the Hickernell writeup.

As before, assume that the function which produced the data  $\mathbf{y}$  is  $f \in \mathcal{H}_K$ . Let us expand the notation for our interpolant  $s$  to include the data  $\mathbf{z}$  which generated the interpolant:  $s(\cdot; \mathbf{z}) = \mathbf{k}(\cdot)^T \mathbf{K}^{-1} \mathbf{z}$ . Also, recall the Hilbert-space norm of an interpolant (derived in (3.4) in the Appendix) is

$$\|s(\cdot; \mathbf{z})\|_{\mathcal{H}_K} = \sqrt{\mathbf{z}^T \mathbf{K}^{-1} \mathbf{z}}. \quad (2.22)$$

Let  $V(\varepsilon)$  denote the volume of the ellipsoid in  $\mathbb{R}^N$  which contains all  $\mathbf{z}$  such that  $\|s(\cdot; \mathbf{z})\|_{\mathcal{H}_K} \leq \|s(\cdot; \mathbf{y})\|_{\mathcal{H}_K}$ :

$$\begin{aligned} V(\varepsilon) &= \text{volume of ellipsoid } \{z \in \mathbb{R}^N : \|s(\cdot; z)\|_{\mathcal{H}_K}^2 \leq \|s(\cdot; \mathbf{y})\|_{\mathcal{H}_K}^2\} \\ &= \text{volume of ellipsoid } \{z \in \mathbb{R}^N : \mathbf{z}^T \mathbf{K}^{-1} \mathbf{z} \leq \mathbf{y}^T \mathbf{K}^{-1} \mathbf{y}\} \\ &= \vartheta_N \frac{(\mathbf{y}^T \mathbf{K}^{-1} \mathbf{y})^N}{\det \mathbf{K}^{-1}} \\ &= \vartheta_N (\mathbf{y}^T \mathbf{K}^{-1} \mathbf{y})^N \det \mathbf{K} = \vartheta_N \exp(C_{\text{MPLE}}(\varepsilon)), \end{aligned}$$

where  $\vartheta_N$  is the volume of the  $\mathbb{R}^N$  unit sphere. Thus, choosing  $\varepsilon$  to minimize the volume of the ellipsoid containing function data which would produce “smaller interpolants” (in the  $\mathcal{H}_K$  norm) than the observed data  $\mathbf{y}$  produces is equivalent to maximizing the profile likelihood that  $\varepsilon$  parameterized the Gaussian field from which  $\mathbf{y}$  was realized.

This concept of ellipsoid volume is basically employing Occam’s Razor: the interpolant that best fits the data should be the simplest, which in this case is measured using the  $\mathcal{H}_K$  norm. The ellipsoid described in  $V(\varepsilon)$  contains data that would produce a simpler interpolant. By choosing  $\varepsilon$  to minimize  $V(\varepsilon)$  we are minimizing the region from which simpler interpolants could be produced, thus making it less likely that our interpolant is not the simplest.

### 2.3.3 Impact of the Hilbert-Schmidt SVD

Two pieces appear in (2.21) which may be subject to ill-conditioning,  $\log(\mathbf{y}^T \mathbf{K}^{-1} \mathbf{y})$  and  $\log \det \mathbf{K}$ . We have already presented (2.8) which will allow the former term to be computed stably, and (2.9) which will bound that term as  $\varepsilon \rightarrow 0$  and help us avoid potential overflow/underflow prior to applying the logarithm. The term  $\log \det \mathbf{K}$  is actually more safe to compute when the analytic form of the eigenvalues is known.

Using the Hilbert-Schmidt SVD,

$$\log \det \mathbf{K} = \log \det \Psi \Lambda_1 \Phi_1^T = \log(\det \Psi \det \Lambda_1 \det \Phi_1^T) = \log \det \Psi + \log \det \Lambda_1 + \log \det \Phi_1^T.$$

Presumably, we have already computed a factorization of  $\Phi_1^T$  when  $\Psi$  was computed (recall (1.9)) so computing  $\log \det \Phi_1^T$  should come at no cost and with no stability issues. Furthermore, assuming we at some point computed  $\Psi^{-1} \mathbf{y}$  to evaluate our interpolant, we should have a factorization of  $\Psi$  meaning that  $\log \det \Psi$  can be computed at no cost and with no stability issues. The final term,  $\log \det \Lambda_1$  is straightforward because it is diagonal:

$$\log \det \Lambda_1 = \log \prod_{n=1}^N \lambda_n = \sum_{n=1}^N \log \lambda_n,$$

and by distributing the logarithm we are able to avoid the underflow issues that would otherwise arise.

[MJM] Add pictures

## 3 Goals for a Parameterization Judgment Tool

For applications to buy into kernel methods, the presence of a free parameter such as  $\varepsilon$  (or perhaps a  $\beta$  smoothness parameter) must be seen as a benefit and not a liability. To do this, we need to have successful parameterization methods so that  $\varepsilon$  can be chosen to produce the most accurate results. The techniques described above show that existing parameterization methods may work or may not work depending on the situation. New schemes can be developed, but we believe that a tool is needed to judge the viability of existing and future schemes within a rigorous context.

Our goal right now is to develop a metric by which the quality of a parameterization scheme can be measured with respect to other schemes, or a single scheme’s success can be compared across scattered data problems. This could be thought of analogously to how we measure the convergence rate of numerical algorithms such as Newton’s method: when solving  $f(\alpha) = 0$ , we know

$$|\epsilon_{n+1}| = \frac{|f''(\xi_n)|}{2|f'(x_n)|} \epsilon_n^2, \quad \epsilon_n = \alpha - x_n, \quad \xi_n \text{ between } x_n \text{ and } \alpha.$$

This suggests that Newton’s method converges quadratically, which is a property that can be compared to other root-finding schemes. On the other hand, this quadratic convergence is only valid when  $f'(x) \neq 0$  near  $x = \alpha$  and  $f''(x) < \infty$  near  $x = \alpha$ , which is a comparison of this method across the set of problems to which it could be applied.

What I think is the likeliest structure for a parameterization judgment tool will be: how near is the guess to the “true” value as the amount of available data increases. Now, what exactly the “true” value means is up for interpretation ...

- You could argue that the true  $\varepsilon$  for any given set of data would be the value that minimizes (some norm of) the error. I think that is an unstable definition because it may be subject to small changes in the data, but it is probably the definition most sought by applications people.
- The “true”  $\varepsilon$  could be the  $\varepsilon$  value that defines the Hilbert-space  $\mathcal{H}_K$  in which  $f$  lies.
- In the same way as the previous bullet, the “true”  $\varepsilon$  could be the  $\varepsilon$  value defining the covariance kernel of the Gaussian process that generated the available data.

The last two points are similar approaches to defining the “true”  $\varepsilon$ , only the 2nd bullet considers it from the numerical analysis side, and the 3rd considers it from the statistics side.

Defining the *fill-distance* to be

$$h_{\mathcal{X}} = \sup_{\mathbf{x} \in \Omega} \min_{\mathbf{x}_j \in \mathcal{X}} \|\mathbf{x} - \mathbf{x}_j\|_2$$

gives us a mechanism for defining the density of the design of points used in the approximation problem. Basically, the fill distance is the radius of the largest ball that can be squeezed in between points in the design  $\mathcal{X}$ . This is a useful tool for studying convergence behavior of interpolation schemes; for instance, we know that a kernel  $K$  with  $2\beta$  smooth derivatives has an interpolant  $s$  to data generated by  $f \in \mathcal{H}_K$  with accuracy

$$|f(\mathbf{x}) - s(\mathbf{x})| \leq Ch_{\mathcal{X}}^{\beta} \sqrt{C_K(\mathbf{x})} \|f\|_{\mathcal{H}_K}, \quad (3.1)$$

where  $C_K$  is independent of  $f$  and  $h_{\mathcal{X}}$  is sufficiently small. An example of an optimal  $\varepsilon$  graph as  $N$  increases (and thus  $h_{\mathcal{X}}$  decreases) is provided in Figure 6.

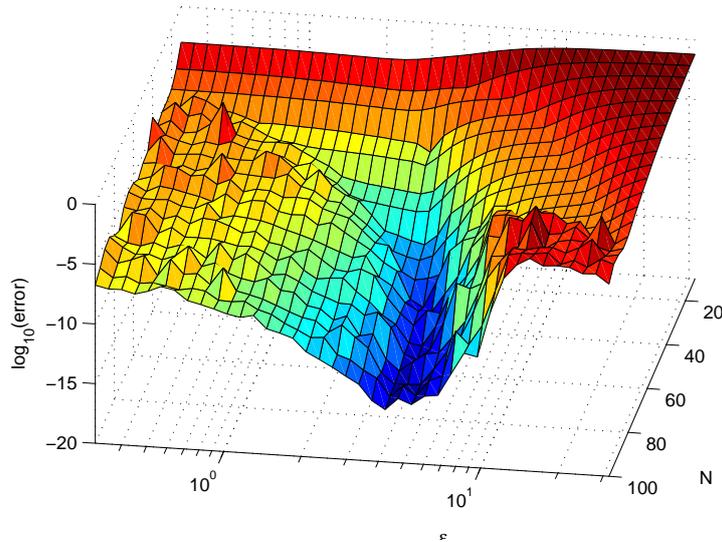


Figure 6: This example involves  $f(x) = (1 + 4x^2)^{-1}$  and shows that as  $N \rightarrow \infty$  there is a rough “convergence” of  $\varepsilon$  to a consistent value. This example doesn’t perfectly reflect the idea of an  $\varepsilon$  value of an underlying RKHS because  $f$  does not belong to a Gaussian RKHS. Even so, the idea still applies here, until machine precision takes over.

It is my hope that a similarly structured bound can exist in a probabilistic sense for the accuracy of a parameterization scheme. Essentially, the situation would play out as

1. Data is generated by a function  $f \in \mathcal{H}_K$  or a Gaussian process  $Y$  with covariance kernel  $K$ .
  - (a) The kernel  $K$  that appears in both settings is the same. It has some parameter  $\varepsilon$  that defines it, but  $\varepsilon$  is unknown to us.

- (b) The data is  $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}$ . This defines the design  $\mathcal{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$  and the vector  $\mathbf{y}$ .
- i. Knowing  $\mathcal{X}$  defines the fill-distance  $h_{\mathcal{X}}$ .
2. A parameterization scheme is chosen to guess a  $\hat{\varepsilon}$  value to be used while constructing the approximation. For example, this could be computed as

$$\hat{\varepsilon} = \underset{\varepsilon}{\operatorname{argmin}} \operatorname{C}_{\text{POWER}}(\varepsilon; 2).$$

- (a) This scheme may involve  $\mathbf{y}$  (e.g., cross-validation) or only  $\mathcal{X}$  (e.g., power function).
- (b) The accuracy of this scheme is hopefully something of the form

$$|\varepsilon - \hat{\varepsilon}| \leq h_{\mathcal{X}}^{\gamma} [C_1(\mathcal{X})C_2(\mathbf{y})\|f\|_{\mathcal{H}_K}],$$

for some  $\gamma > 0$  and  $C_1, C_2$ . Really, I have no idea if this is what it will look like, it's just a thought for an ideal situation in the deterministic setting.

- (c) What is much more likely is that the bound will be of the form

$$P(|\varepsilon - \hat{\varepsilon}| < \alpha) \geq 1 - \nu(\mathcal{X}, \mathbf{y}, \alpha), \quad \lim_{h_{\mathcal{X}} \rightarrow 0} \nu(\mathcal{X}, \mathbf{y}, \alpha) = 0,$$

because if the data is generated by a Gaussian process, then there's always a chance that it will be a really crummy, uninformative realization.

This last point, Dr. Cobb, is where I think you will be of great help. If we can prove something of this form, then we would be able to judge (empirically at the least) the quality of parameterization schemes which would help inform applications scientists.

With a tool like this, we would be able to make comments on

- **Consistency** - Will the scheme recover the “true”  $\varepsilon$  for an infinitely dense design? This would only be true if  $\lim_{h_{\mathcal{X}} \rightarrow 0} \nu(\mathcal{X}, \mathbf{y}) = 0$ .
- **Convergence rate** - How quickly is the parameterization scheme approaching the “true”  $\varepsilon$ , and thus how few points are needed before I feel comfortable that I am doing a decent job? This would depend on the rate at which  $\lim_{h_{\mathcal{X}} \rightarrow 0} \nu(\mathcal{X}, \mathbf{y}) = 0$ , especially in comparison to other schemes.
- **Stability** - Do small changes in  $\mathcal{X}$  and  $\mathbf{y}$  affect the consistency or convergence of the scheme? I'm not as concerned about this, but it is something that should eventually be studied.
- **Bounding  $\varepsilon$**  - Most applications don't demand an optimal  $\varepsilon$  because of noise in the data. Can this convergence study be used to create a region in which the “true”  $\varepsilon$  lies?
- **Computational cost** - If two schemes have similar convergence properties, is there any reason to not use the cheaper one?
- **Log scale** - The plots used above show  $\varepsilon$  on a log-scale, and often this is how we consider different  $\varepsilon$  values. Is it useful to instead study accuracy of the form

$$|\log \varepsilon - \log \hat{\varepsilon}|, \quad \text{or} \quad \log |\varepsilon - \hat{\varepsilon}|?$$

So, do you know of anything that would help us think about this problem?

## References

- [1] R. Cavoretto, G. Fasshauer, and M. McCourt. Compact Matérn kernels and piecewise polynomial splines viewed from a Hilbert-Schmidt perspective, 2013. submitted.
- [2] G. E. Fasshauer. *Meshfree Approximation Methods with MATLAB*. World Scientific Publishing Co., Inc., River Edge, NJ, USA, 2007.

- [3] G. E. Fasshauer and M. McCourt. Stable evaluation of Gaussian RBF interpolants. *SIAM J. Sci. Comput.*, 34(2):A737–A762, 2012.
- [4] M. Golomb and H. F. Weinberger. Optimal approximation and error bounds. In R. E. Langer, editor, *On Numerical Approximation*, pages 117–190. University of Wisconsin Press, 1959.
- [5] G. H. Golub and C. F. Van Loan. *Matrix Computations*. Johns Hopkins Studies in the Mathematical Sciences. Johns Hopkins University Press, 2013.
- [6] M. McCourt. Using Gaussian eigenfunctions to solve boundary value problems. *Adv. Appl. Math. Mech.*, 5(4):569–594, 2013.
- [7] M. L. Stein. *Interpolation of Spatial Data: Some Theory for Kriging*. Springer Series in Statistics. Springer, 1999.
- [8] H. Wendland. *Scattered Data Approximation*. Cambridge Monographs on Applied and Computational Mathematics. Cambridge University Press, 2005.

## Appendix: Deriving the Native Space Error Bound

For this derivation, we need to recall the reproducing property

$$\langle f, K(\cdot, \mathbf{x}) \rangle_{\mathcal{H}_K} = f(\mathbf{x}), \quad f \in \mathcal{H}_K,$$

where  $\mathcal{H}_K$  is the reproducing kernel Hilbert space, or *native space*, generated by  $K$ . It is important to also know that  $K(\cdot, x) \in \mathcal{H}_K$  for  $x \in \Omega$ , which implies that

$$\langle K(\cdot, \mathbf{x}), K(\cdot, \mathbf{z}) \rangle_{\mathcal{H}_K} = K(\mathbf{x}, \mathbf{z}), \quad \mathbf{x}, \mathbf{z} \in \Omega.$$

Note that because  $s(\mathbf{x})$  is a number,

$$s(\mathbf{x}) = s(\mathbf{x})^T = \mathbf{y}^T \mathbf{K}^{-1} \mathbf{k}(\mathbf{x}).$$

The vector  $\mathbf{y}$  is full of function values ( $y_i = f(\mathbf{x}_i)$ ) so we can write

$$\begin{aligned} \mathbf{y}^T &= (f(\mathbf{x}_1) \quad \cdots \quad f(\mathbf{x}_N)) = (\langle f, K(\cdot, \mathbf{x}_1) \rangle_{\mathcal{H}_K} \quad \cdots \quad \langle f, K(\cdot, \mathbf{x}_N) \rangle_{\mathcal{H}_K}) \\ &= \langle f, (K(\cdot, \mathbf{x}_1) \quad \cdots \quad K(\cdot, \mathbf{x}_N)) \rangle_{\mathcal{H}_K} \\ &= \langle f, \mathbf{k}(\cdot)^T \rangle_{\mathcal{H}_K}. \end{aligned}$$

That in turn allows us to write

$$s(\mathbf{x}) = \langle f, \mathbf{k}(\cdot)^T \rangle_{\mathcal{H}_K} \mathbf{K}^{-1} \mathbf{k}(\mathbf{x}) = \langle f, \mathbf{k}(\cdot)^T \mathbf{K}^{-1} \mathbf{k}(\mathbf{x}) \rangle_{\mathcal{H}_K}.$$

because  $\mathbf{K}^{-1} \mathbf{k}(\mathbf{x})$  is not a function and is unaffected by the inner product. Using this in the error expression gives

$$\begin{aligned} |f(\mathbf{x}) - s(\mathbf{x})| &= |f(\mathbf{x}) - \mathbf{y}^T \mathbf{K}^{-1} \mathbf{k}(\mathbf{x})| \\ &= \left| \langle f, K(\cdot, \mathbf{x}) \rangle_{\mathcal{H}_K} - \langle f, \mathbf{k}(\cdot)^T \mathbf{K}^{-1} \mathbf{k}(\mathbf{x}) \rangle_{\mathcal{H}_K} \right| \\ &= \left| \langle f, K(\cdot, \mathbf{x}) - \mathbf{k}(\cdot)^T \mathbf{K}^{-1} \mathbf{k}(\mathbf{x}) \rangle_{\mathcal{H}_K} \right| \\ &\leq \|f\|_{\mathcal{H}_K} \left\| K(\cdot, \mathbf{x}) - \mathbf{k}^T(\cdot) \mathbf{K}^{-1} \mathbf{k}(\mathbf{x}) \right\|_{\mathcal{H}_K} = \|f\|_{\mathcal{H}_K} P_{K, \mathcal{X}}(\mathbf{x}), \end{aligned}$$

using the Cauchy-Schwarz inequality, and the power function

$$P_{K, \mathcal{X}}(\mathbf{x}) = \sqrt{K(\mathbf{x}, \mathbf{x}) - \mathbf{k}^T(\mathbf{x}) \mathbf{K}^{-1} \mathbf{k}(\mathbf{x})}. \quad (1.6)$$

This expression for the power function is obtained by again using the reproducing property on the norm present in the inequality.

The standard error bound derived above

$$|f(\mathbf{x}) - s(\mathbf{x})| \leq \|f\|_{\mathcal{H}_K} P_{K,\mathcal{X}}(\mathbf{x}) \quad (1.5)$$

can be improved (see [4]) to

$$|f(\mathbf{x}) - s(\mathbf{x})| \leq \|f - s\|_{\mathcal{H}_K} P_{K,\mathcal{X}}(\mathbf{x}). \quad (3.2)$$

To see this, first we must prove that  $f - s$  is orthogonal to  $s$  in the Hilbert-space inner product:

$$\begin{aligned} \langle f - s, s \rangle_{\mathcal{H}_K} &= \langle f - s, \mathbf{k}(\cdot)^T \mathbf{K}^{-1} \mathbf{y} \rangle_{\mathcal{H}_K} \\ &= \langle f - s, \mathbf{k}(\cdot)^T \rangle_{\mathcal{H}_K} \mathbf{K}^{-1} \mathbf{y} \\ &= (\langle f - s, K(\cdot, \mathbf{x}_1) \rangle_{\mathcal{H}_K} \cdots \langle f - s, K(\cdot, \mathbf{x}_N) \rangle_{\mathcal{H}_K}) \mathbf{K}^{-1} \mathbf{y} \\ &= (f(\mathbf{x}_1) - s(\mathbf{x}_1)) \cdots (f(\mathbf{x}_N) - s(\mathbf{x}_N)) \mathbf{K}^{-1} \mathbf{y} = (0 \cdots 0) \mathbf{K}^{-1} \mathbf{y} = 0 \end{aligned}$$

because we know that  $f(\mathbf{x}_i) = s(\mathbf{x}_i)$ ,  $1 \leq i \leq N$ , if  $s$  interpolates  $f$  at those points. Using this, we can prove that  $\|f - s\|_{\mathcal{H}_K} < \|f\|_{\mathcal{H}_K}$ :

$$\|f\|_{\mathcal{H}_K}^2 = \|f - s + s\|_{\mathcal{H}_K}^2 = \|f - s\|_{\mathcal{H}_K}^2 + 2\langle f - s, s \rangle_{\mathcal{H}_K} + \|s\|_{\mathcal{H}_K}^2 = \|f - s\|_{\mathcal{H}_K}^2 + \|s\|_{\mathcal{H}_K}^2 > \|f - s\|_{\mathcal{H}_K}^2, \quad (3.3)$$

which implies that (3.2) is a better bound. Even so, this tighter error bound does not seem to play a significant role in the kernel literature.

Since  $\|f\|_{\mathcal{H}(K,\Omega)}$  usually is not computable (remember, we do not even know  $f$ , but want to reconstruct it from the data) the standard error bound is not very useful for practical situations. On the other hand, if we assume that our approximation  $s$  is good, i.e.,

$$\|f - s\|_{\mathcal{H}_K} \leq \delta_\varepsilon \|s\|_{\mathcal{H}_K}$$

for some not too large constant  $\delta_\varepsilon$ , then the Golomb-Weinberger improved error bound yields a mostly computable error bound

$$|f(\mathbf{x}) - s(\mathbf{x})| \leq \delta_\varepsilon \|s\|_{\mathcal{H}_K} P_{K,\mathcal{X}}(\mathbf{x}).$$

This is indeed computable since  $\|s\|_{\mathcal{H}_K} = \sqrt{\mathbf{y}^T \mathbf{K}^{-1} \mathbf{y}}$ :

$$\begin{aligned} \|s\|_{\mathcal{H}_K}^2 &= \langle s, s \rangle_{\mathcal{H}_K} = \langle \mathbf{y}^T \mathbf{K}^{-1} \mathbf{k}(\cdot), \mathbf{k}(\cdot)^T \mathbf{K}^{-1} \mathbf{y} \rangle_{\mathcal{H}_K} \\ &= \mathbf{y}^T \mathbf{K}^{-1} \langle \mathbf{k}(\cdot), \mathbf{k}(\cdot)^T \rangle_{\mathcal{H}_K} \mathbf{K}^{-1} \mathbf{y} = \mathbf{y}^T \mathbf{K}^{-1} \mathbf{K} \mathbf{K}^{-1} \mathbf{y} = \mathbf{y}^T \mathbf{K}^{-1} \mathbf{y}. \end{aligned} \quad (3.4)$$

Also, recall that this term  $\mathbf{y}^T \mathbf{K}^{-1} \mathbf{y}$  appears in (1.17) for the zero-mean Gaussian field.