

Some Recent Insights into Computing with Positive Definite Kernels

Greg Fasshauer

Department of Applied Mathematics & Statistics
Colorado School of Mines
Partially supported by NSF Grant DMS-1522687

Joint work with Mike McCourt

MAIA Luminy, September 21, 2016



Outline

- 1 Deterministic and Statistical Views of Kernel Methods
- 2 Parametrization Criteria
- 3 Computational Aspects
- 4 Numerical Illustrations

Outline

- 1 Deterministic and Statistical Views of Kernel Methods
- 2 Parametrization Criteria
- 3 Computational Aspects
- 4 Numerical Illustrations

New Perspectives



Illinois

New Perspectives



Colorado

Deterministic Kernel-based Interpolation

Given data $\{(\mathbf{x}_i, y_i)\}_{i=1}^N$, use a data-dependent linear function space, i.e.,

$$s(\mathbf{x}) = \sum_{j=1}^N c_j K(\mathbf{x}, \mathbf{x}_j) = \mathbf{k}(\mathbf{x})^T \mathbf{c}, \quad \mathbf{x} \in \Omega \subseteq \mathbb{R}^d$$

with $K : \Omega \times \Omega \rightarrow \mathbb{R}$ a **positive definite reproducing kernel**.

Deterministic Kernel-based Interpolation

Given data $\{(\mathbf{x}_i, y_i)\}_{i=1}^N$, use a data-dependent linear function space, i.e.,

$$\mathbf{s}(\mathbf{x}) = \sum_{j=1}^N c_j K(\mathbf{x}, \mathbf{x}_j) = \mathbf{k}(\mathbf{x})^T \mathbf{c}, \quad \mathbf{x} \in \Omega \subseteq \mathbb{R}^d$$

with $K : \Omega \times \Omega \rightarrow \mathbb{R}$ a **positive definite reproducing kernel**.

To find c_j solve the interpolation equations

$$\mathbf{s}(\mathbf{x}_i) = y_i, \quad i = 1, \dots, N,$$

which leads to a linear system $\mathbf{K}\mathbf{c} = \mathbf{y}$ with symmetric positive definite — **often ill-conditioned** — system matrix

$$K_{ij} = K(\mathbf{x}_i, \mathbf{x}_j), \quad i, j = 1, \dots, N,$$

Deterministic Kernel-based Interpolation

Given data $\{(\mathbf{x}_i, y_i)\}_{i=1}^N$, use a data-dependent linear function space, i.e.,

$$s(\mathbf{x}) = \sum_{j=1}^N c_j K(\mathbf{x}, \mathbf{x}_j) = \mathbf{k}(\mathbf{x})^T \mathbf{c}, \quad \mathbf{x} \in \Omega \subseteq \mathbb{R}^d$$

with $K : \Omega \times \Omega \rightarrow \mathbb{R}$ a **positive definite reproducing kernel**.

To find c_j solve the interpolation equations

$$s(\mathbf{x}_i) = y_i, \quad i = 1, \dots, N,$$

which leads to a linear system $\mathbf{K}\mathbf{c} = \mathbf{y}$ with symmetric positive definite — **often ill-conditioned** — system matrix

$$K_{ij} = K(\mathbf{x}_i, \mathbf{x}_j), \quad i, j = 1, \dots, N,$$

and **cardinal representation**

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

Connection to Kriging

Given a Gaussian (zero-mean) random field Y with covariance kernel K , process variance σ^2 and observations

$$\mathbf{Y} = (Y_{\mathbf{x}_1} \quad \cdots \quad Y_{\mathbf{x}_N})^T, \quad Y_{\mathbf{x}_j} \text{ zero-mean random variables,}$$

the (simple) **kriging predictor** is of the form

$$\hat{Y}_{\mathbf{x}} = \sum_{j=1}^N w_j(\mathbf{x}) Y_{\mathbf{x}_j} = \mathbf{w}(\mathbf{x})^T \mathbf{Y},$$

- $\hat{Y}_{\mathbf{x}}$: **zero-mean random variable**,
- $\mathbf{w}(\cdot) = (w_1(\cdot) \quad \cdots \quad w_N(\cdot))^T$: vector of **weight functions**.

“Optimal” weights $w_j^*(\cdot)$ will minimize the MSE of the predictor, i.e.,

$$\text{MSE}(\hat{Y}_{\mathbf{x}}) = \mathbb{E} \left[\left(Y_{\mathbf{x}} - \mathbf{w}(\mathbf{x})^T \mathbf{Y} \right)^2 \right].$$



Using the covariance kernel K and process variance σ^2 of Y , i.e., $\sigma^2 K(\mathbf{x}, \mathbf{z}) = \mathbb{E}[Y_{\mathbf{x}} Y_{\mathbf{z}}]$, we have

$$\begin{aligned} \text{MSE}(\hat{Y}_{\mathbf{x}}) &= \mathbb{E} \left[\left(Y_{\mathbf{x}} - \mathbf{w}(\mathbf{x})^T \mathbf{Y} \right)^2 \right] \\ &= \mathbb{E}[Y_{\mathbf{x}} Y_{\mathbf{x}}] - 2\mathbb{E}[Y_{\mathbf{x}} \mathbf{w}(\mathbf{x})^T \mathbf{Y}] + \mathbb{E}[\mathbf{w}(\mathbf{x})^T \mathbf{Y} \mathbf{Y}^T \mathbf{w}(\mathbf{x})] \\ &= \sigma^2 K(\mathbf{x}, \mathbf{x}) - 2\mathbf{w}(\mathbf{x})^T (\sigma^2 \mathbf{k}(\mathbf{x})) + \mathbf{w}(\mathbf{x})^T (\sigma^2 \mathbf{K}) \mathbf{w}(\mathbf{x}). \end{aligned}$$

Using the covariance kernel K and process variance σ^2 of Y , i.e., $\sigma^2 K(\mathbf{x}, \mathbf{z}) = \mathbb{E}[Y_{\mathbf{x}} Y_{\mathbf{z}}]$, we have

$$\begin{aligned} \text{MSE}(\hat{Y}_{\mathbf{x}}) &= \mathbb{E} \left[\left(Y_{\mathbf{x}} - \mathbf{w}(\mathbf{x})^T \mathbf{Y} \right)^2 \right] \\ &= \mathbb{E}[Y_{\mathbf{x}} Y_{\mathbf{x}}] - 2\mathbb{E}[Y_{\mathbf{x}} \mathbf{w}(\mathbf{x})^T \mathbf{Y}] + \mathbb{E}[\mathbf{w}(\mathbf{x})^T \mathbf{Y} \mathbf{Y}^T \mathbf{w}(\mathbf{x})] \\ &= \sigma^2 K(\mathbf{x}, \mathbf{x}) - 2\mathbf{w}(\mathbf{x})^T (\sigma^2 \mathbf{k}(\mathbf{x})) + \mathbf{w}(\mathbf{x})^T (\sigma^2 \mathbf{K}) \mathbf{w}(\mathbf{x}). \end{aligned}$$

Differentiation and equating to 0 yields the optimum weight vector

$$\hat{\mathbf{w}}(\mathbf{x}) = \mathbf{K}^{-1} \mathbf{k}(\mathbf{x}),$$

Using the covariance kernel K and process variance σ^2 of Y , i.e., $\sigma^2 K(\mathbf{x}, \mathbf{z}) = \mathbb{E}[Y_{\mathbf{x}} Y_{\mathbf{z}}]$, we have

$$\begin{aligned} \text{MSE}(\hat{Y}_{\mathbf{x}}) &= \mathbb{E} \left[\left(Y_{\mathbf{x}} - \mathbf{w}(\mathbf{x})^T \mathbf{Y} \right)^2 \right] \\ &= \mathbb{E}[Y_{\mathbf{x}} Y_{\mathbf{x}}] - 2\mathbb{E}[Y_{\mathbf{x}} \mathbf{w}(\mathbf{x})^T \mathbf{Y}] + \mathbb{E}[\mathbf{w}(\mathbf{x})^T \mathbf{Y} \mathbf{Y}^T \mathbf{w}(\mathbf{x})] \\ &= \sigma^2 K(\mathbf{x}, \mathbf{x}) - 2\mathbf{w}(\mathbf{x})^T (\sigma^2 \mathbf{k}(\mathbf{x})) + \mathbf{w}(\mathbf{x})^T (\sigma^2 \mathbf{K}) \mathbf{w}(\mathbf{x}). \end{aligned}$$

Differentiation and equating to 0 yields the optimum weight vector

$$\hat{\mathbf{w}}(\mathbf{x}) = \mathbf{K}^{-1} \mathbf{k}(\mathbf{x}),$$

so that the (simple) kriging predictor

$$\hat{Y}_{\mathbf{x}} = \mathbf{k}(\mathbf{x})^T \mathbf{K}^{-1} \mathbf{Y}$$

is the best linear unbiased predictor.

Remark

Note that *this is independent of the process variance σ^2 .*

Parametrized Kernel Methods

Many kernels contain **parameters whose values greatly affect the performance of the kernel method.**

For example, such parameters may affect

- shape,
- smoothness,
- accuracy,
- numerical stability,
- computational efficiency, i.e., density of (sparse) K ,
- balance between closeness of fit and smoothness.

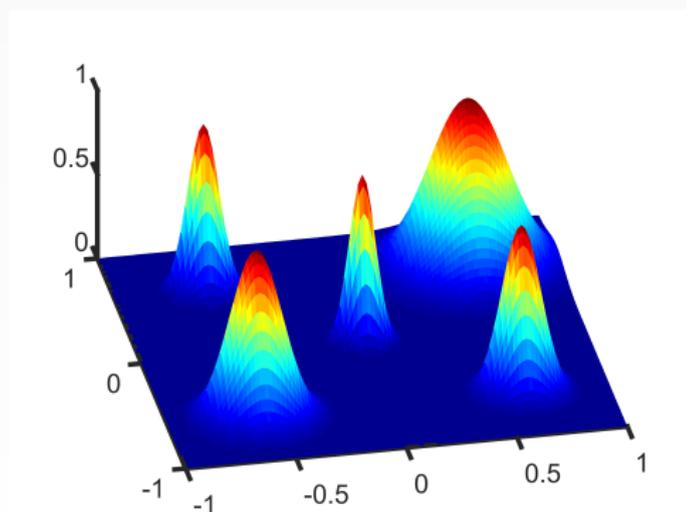
It is therefore important to **find “optimal” values for these parameters.**



Examples

- isotropic shape parameter, e.g.,

$$K(\mathbf{x}, \mathbf{z}) = \kappa(r) = e^{-\varepsilon^2 r^2}, \quad r = \|\mathbf{x} - \mathbf{z}\|_2$$

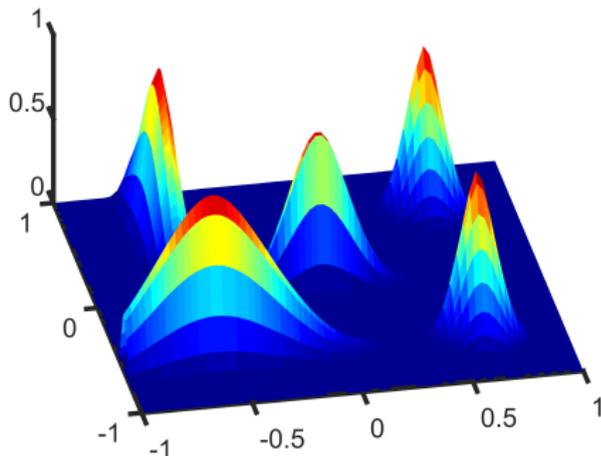


Isotropic Gaussian kernels, $\varepsilon = 4, 6, 8, 10, 12$

Examples (cont.)

- radial anisotropic shape parameters, e.g.,

$$\kappa(r) = (1-r)_+^4(4r+1), \quad r = \sqrt{(\mathbf{x} - \mathbf{z})^T \mathbf{E} (\mathbf{x} - \mathbf{z})}, \quad \mathbf{E} = \begin{pmatrix} \varepsilon_1^2 & & \\ & \ddots & \\ & & \varepsilon_d^2 \end{pmatrix}$$

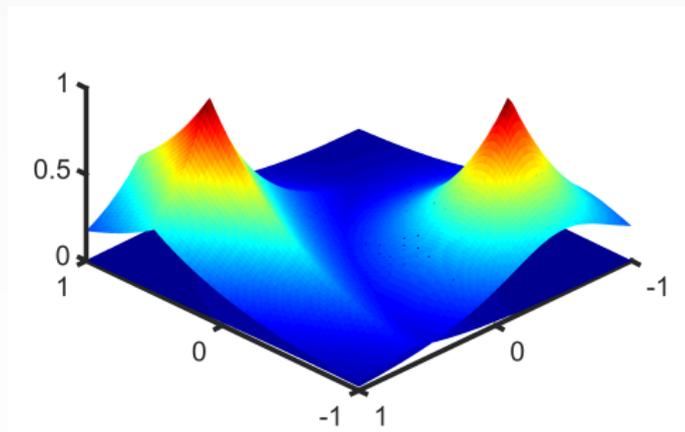


Compactly supported Wendland kernels,
 $\varepsilon = [4, 4], [2, 2], [5, 1], [1, 5], [2, 7]$

Examples (cont.)

- tensor product anisotropic shape parameters, e.g.,

$$K(\mathbf{x}, \mathbf{z}) = \prod_{\ell=1}^d e^{-\varepsilon_{\ell} |x_{\ell} - z_{\ell}|}$$



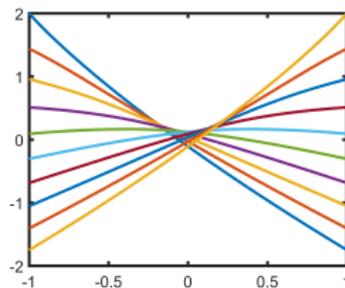
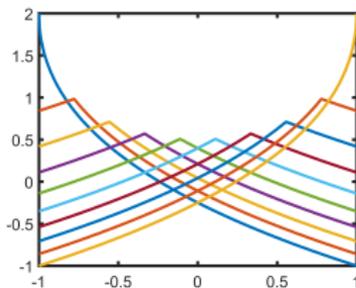
Tensor product and radial C^0 Matérn kernels,
 $\varepsilon = [1.5, 2.5], \varepsilon = 2.5$

Examples (cont.)

- smoothness parameter(s), e.g.,

$$K(x, z) = \sum_{n=1}^{\infty} \lambda_n \varphi_n(x) \varphi_n(z)$$

$$\lambda_n = \frac{1}{\zeta(2\beta)n^{2\beta}}, \quad \varphi_n(x) = \sqrt{2}T_n(x), \quad n = 1, 2, \dots$$



Chebyshev kernels with $\beta = 1, 2$

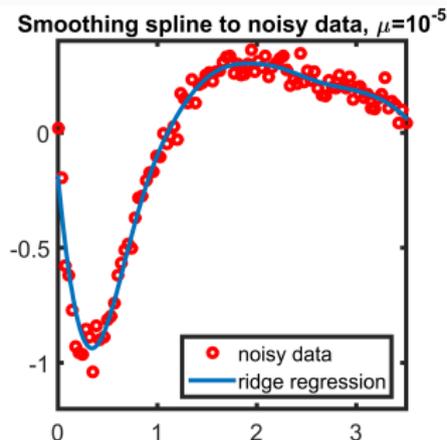
Examples (cont.)

- regularization parameter, e.g.,

$$\min_{\mathbf{c} \in \mathbb{R}^N} \left[(\mathbf{y} - \mathbf{K}\mathbf{c})^T (\mathbf{y} - \mathbf{K}\mathbf{c}) + \mu \mathbf{c}^T \mathbf{K}\mathbf{c} \right]$$

so that we have the **smoothing spline/ridge regression**

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



Examples (cont.)

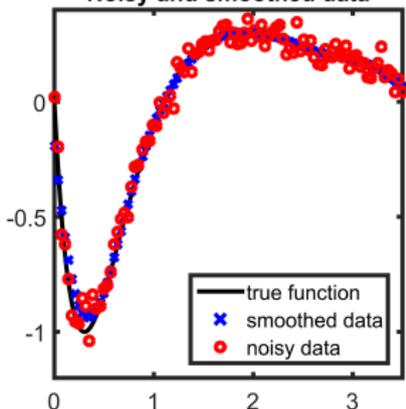
- regularization parameter, e.g.,

$$\min_{\mathbf{c} \in \mathbb{R}^N} \left[(\mathbf{y} - \mathbf{K}\mathbf{c})^T (\mathbf{y} - \mathbf{K}\mathbf{c}) + \mu \mathbf{c}^T \mathbf{K} \mathbf{c} \right]$$

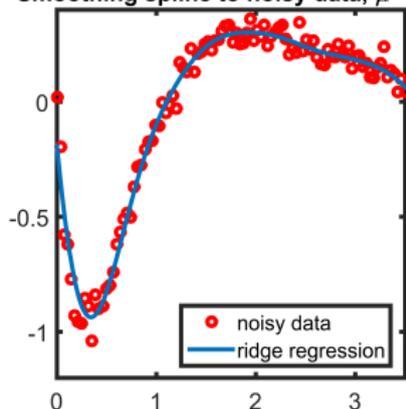
so that we have the **smoothing spline/ridge regression**

$$\begin{aligned} s(\mathbf{x}) &= \mathbf{k}(\mathbf{x})^T (\mathbf{K} + \mu \mathbf{I})^{-1} \mathbf{y} = \mathbf{k}(\mathbf{x})^T \mathbf{K}^{-1} (\mathbf{K} + \mu \mathbf{I})^{-1} \mathbf{K} \mathbf{y} \\ &= \ell(\mathbf{x})^T \tilde{\mathbf{y}}. \end{aligned}$$

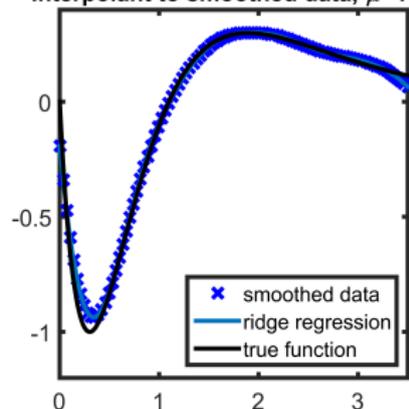
Noisy and smoothed data



Smoothing spline to noisy data, $\mu=10^{-5}$



Interpolant to smoothed data, $\mu=10^{-5}$



Outline

- 1 Deterministic and Statistical Views of Kernel Methods
- 2 Parametrization Criteria**
- 3 Computational Aspects
- 4 Numerical Illustrations

How do we decide “optimality”?

We use parametrization criteria such as

$$C_{\text{CV}}(\varepsilon; \rho) = \left\| \left(\frac{\mathbf{c}_1}{K_{11}^{-1}} \quad \cdots \quad \frac{\mathbf{c}_N}{K_{NN}^{-1}} \right) \right\|_{\rho},$$

LOOCV

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

profile likelihood

$$C_{\text{GW}}(\varepsilon; \rho) = \sqrt{\mathbf{y}^T \mathbf{K}^{-1} \mathbf{y}} \|P_{K, \mathbf{x}}\|_{\rho},$$

Golomb–Weinberger

where $P_{K, \mathbf{x}}(\mathbf{x}) = \sqrt{K(\mathbf{x}, \mathbf{x}) - \mathbf{k}(\mathbf{x})^T \mathbf{K}^{-1} \mathbf{k}(\mathbf{x})}$

or kriging variance

These criteria all aim to maximize some form of accuracy.

They may require computing

$$\mathbf{K}^{-1} \mathbf{y}, \quad \mathbf{y}^T \mathbf{K}^{-1} \mathbf{y}, \quad \mathbf{k}(\mathbf{x})^T \mathbf{K}^{-1}, \quad \log \det \mathbf{K}, \quad P_{K, \mathbf{x}}(\mathbf{x}).$$



Outline for remainder of talk

- Explain parametrization criteria:
 - maximum profile likelihood
 - Golomb–Weinberger error bound (kriging variance)
- Explain how to compute
 - cardinal functions $\ell(\mathbf{x})^T = \mathbf{k}(\mathbf{x})^T \mathbf{K}^{-1}$
 - log det \mathbf{K}
 - power function $P_{\mathbf{K}, \mathbf{x}}(\mathbf{x}) = \sqrt{\mathbf{K}(\mathbf{x}, \mathbf{x}) - \mathbf{k}(\mathbf{x})^T \mathbf{K}^{-1} \mathbf{k}(\mathbf{x})}$
- Show a few numerical examples

Maximum Profile Likelihood

We consider a zero-mean Gaussian random field Y with covariance kernel $\sigma^2 K$ and likelihood function

$$L(\sigma^2, \varepsilon) = \frac{1}{\sqrt{(2\pi)^N \sigma^{2N} \det(K)}} \exp\left(-\frac{1}{2\sigma^2} \mathbf{y}^T K^{-1} \mathbf{y}\right),$$

where ε appears within K .

Maximum Profile Likelihood

We consider a zero-mean Gaussian random field \mathbf{Y} with covariance kernel $\sigma^2 \mathbf{K}$ and likelihood function

$$L(\sigma^2, \varepsilon) = \frac{1}{\sqrt{(2\pi)^N \sigma^{2N} \det(\mathbf{K})}} \exp\left(-\frac{1}{2\sigma^2} \mathbf{y}^T \mathbf{K}^{-1} \mathbf{y}\right),$$

where ε appears within \mathbf{K} .

Maximizing $L(\sigma^2, \varepsilon)$ is equivalent to minimizing

$$-2 \log\left(L(\sigma^2, \varepsilon)\right) = N \log 2\pi + N \log \sigma^2 + \log \det \mathbf{K} + \frac{1}{\sigma^2} \mathbf{y}^T \mathbf{K}^{-1} \mathbf{y}.$$

Maximum Profile Likelihood

We consider a **zero-mean Gaussian random field** \mathbf{Y} with **covariance kernel** $\sigma^2 \mathbf{K}$ and **likelihood function**

$$L(\sigma^2, \varepsilon) = \frac{1}{\sqrt{(2\pi)^N \sigma^{2N} \det(\mathbf{K})}} \exp\left(-\frac{1}{2\sigma^2} \mathbf{y}^T \mathbf{K}^{-1} \mathbf{y}\right),$$

where ε appears within \mathbf{K} .

Maximizing $L(\sigma^2, \varepsilon)$ is equivalent to **minimizing**

$$-2 \log\left(L(\sigma^2, \varepsilon)\right) = N \log 2\pi + N \log \sigma^2 + \log \det \mathbf{K} + \frac{1}{\sigma^2} \mathbf{y}^T \mathbf{K}^{-1} \mathbf{y}.$$

Differentiating with respect to σ^2 and equating to 0 gives the **optimal profile variance**

$$\sigma_{\text{opt}}^2 = \frac{1}{N} \mathbf{y}^T \mathbf{K}^{-1} \mathbf{y}.$$



Using the **optimal profile variance** $\sigma_{\text{opt}}^2 = \frac{1}{N} \mathbf{y}^T \mathbf{K}^{-1} \mathbf{y}$, we now find ε that minimizes

$$\begin{aligned}
 & -2 \log \left(L(\sigma_{\text{opt}}^2, \varepsilon) \right) - N \log 2\pi \\
 & = N \log \sigma_{\text{opt}}^2 + \log \det \mathbf{K} + \frac{1}{\sigma_{\text{opt}}^2} \mathbf{y}^T \mathbf{K}^{-1} \mathbf{y} \\
 & = N \log \left(\mathbf{y}^T \mathbf{K}^{-1} \mathbf{y} \right) - N \log N + \log \det \mathbf{K} + N.
 \end{aligned}$$

Using the **optimal profile variance** $\sigma_{\text{opt}}^2 = \frac{1}{N} \mathbf{y}^T \mathbf{K}^{-1} \mathbf{y}$, we now find ε that minimizes

$$\begin{aligned} & -2 \log \left(L(\sigma_{\text{opt}}^2, \varepsilon) \right) - N \log 2\pi \\ &= N \log \sigma_{\text{opt}}^2 + \log \det \mathbf{K} + \frac{1}{\sigma_{\text{opt}}^2} \mathbf{y}^T \mathbf{K}^{-1} \mathbf{y} \\ &= N \log \left(\mathbf{y}^T \mathbf{K}^{-1} \mathbf{y} \right) - N \log N + \log \det \mathbf{K} + N. \end{aligned}$$

This yields the **profile likelihood criterion**

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

Golomb–Weinberger/kriging variance criterion

Using the representation $\ell(\mathbf{x})^T = \mathbf{k}(\mathbf{x})^T \mathbf{K}^{-1}$ of the cardinal functions and

$$\mathbf{y} = (f(\mathbf{x}_1) \quad \cdots \quad f(\mathbf{x}_N))^T$$

$$\stackrel{K \text{ is RK}}{=} (\langle f, K(\cdot, \mathbf{x}_1) \rangle_{\mathcal{H}_K} \quad \cdots \quad \langle f, K(\cdot, \mathbf{x}_N) \rangle_{\mathcal{H}_K})^T = \langle f, \mathbf{k}(\cdot) \rangle_{\mathcal{H}_K}$$

we have the **standard pointwise error bound** for interpolation

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

with **power function**

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

The standard error bound

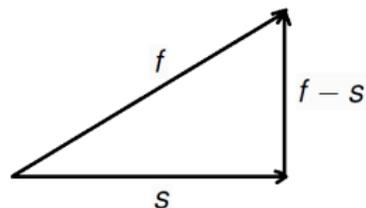
$$|f(\mathbf{x}) - s(\mathbf{x})| \leq \|f\|_{\mathcal{H}_K} P_{K,\mathcal{X}}(\mathbf{x})$$

can be improved to the tight bound (see [GW59])

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

since $f - s$ is orthogonal to s in \mathcal{H}_K , i.e.,

$$\|f\|_{\mathcal{H}_K}^2 = \|f - s\|_{\mathcal{H}_K}^2 + \|s\|_{\mathcal{H}_K}^2.$$



The **standard error bound**

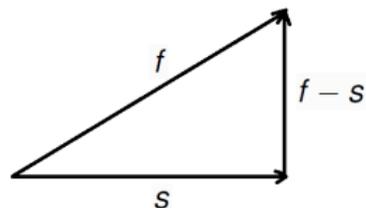
$$|f(\mathbf{x}) - s(\mathbf{x})| \leq \|f\|_{\mathcal{H}_K} P_{K,\mathcal{X}}(\mathbf{x})$$

can be improved to the tight bound (see [GW59])

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

since $f - s$ is orthogonal to s in \mathcal{H}_K , i.e.,

$$\|f\|_{\mathcal{H}_K}^2 = \|f - s\|_{\mathcal{H}_K}^2 + \|s\|_{\mathcal{H}_K}^2.$$



Problem: Usually, **neither** $\|f\|_{\mathcal{H}_K}$ **nor** $\|f - s\|_{\mathcal{H}_K}$ **are computable.**

The **standard error bound**

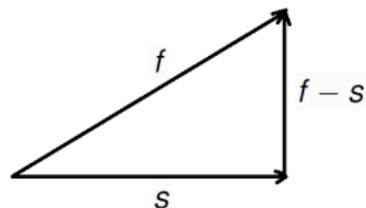
$$|f(\mathbf{x}) - s(\mathbf{x})| \leq \|f\|_{\mathcal{H}_K} P_{K,\mathcal{X}}(\mathbf{x})$$

can be **improved** to the tight bound (see [GW59])

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

since $f - s$ is orthogonal to s in \mathcal{H}_K , i.e.,

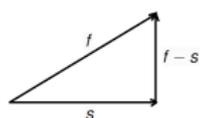
$$\|f\|_{\mathcal{H}_K}^2 = \|f - s\|_{\mathcal{H}_K}^2 + \|s\|_{\mathcal{H}_K}^2.$$



Problem: Usually, **neither** $\|f\|_{\mathcal{H}_K}$ **nor** $\|f - s\|_{\mathcal{H}_K}$ **are computable.**
But $\|s\|_{\mathcal{H}_K}$ **is.**

If we assume that our approximation s is not too bad, i.e.,

$$\|f - s\|_{\mathcal{H}_K(\Omega)} \leq \delta \|s\|_{\mathcal{H}_K(\Omega)}$$

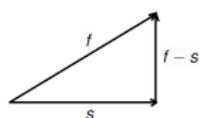


for some not too large constant δ , then the Golomb–Weinberger tight error bound yields a **computable error bound**

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

If we assume that our approximation s is not too bad, i.e.,

$$\|f - s\|_{\mathcal{H}_K(\Omega)} \leq \delta \|s\|_{\mathcal{H}_K(\Omega)}$$



for some not too large constant δ , then the Golomb–Weinberger tight error bound yields a **computable error bound**

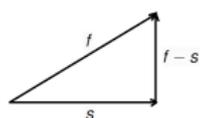
$$|f(\mathbf{x}) - s(\mathbf{x})| \leq \delta \|s\|_{\mathcal{H}_K} P_{K,\mathcal{X}}(\mathbf{x}).$$

This is indeed computable since

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

If we assume that our approximation s is not too bad, i.e.,

$$\|f - s\|_{\mathcal{H}_K(\Omega)} \leq \delta \|s\|_{\mathcal{H}_K(\Omega)}$$



for some not too large constant δ , then the Golomb–Weinberger tight error bound yields a **computable error bound**

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

This is indeed computable since

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

Therefore we have

$$C_{\text{GW}}(\varepsilon; \rho) = \sqrt{\mathbf{y}^T \mathbf{K}^{-1} \mathbf{y}} \|P_{K,\mathcal{X}}\|_p,$$

where we **compute the p -norm on a discrete evaluation grid.**

Kriging Variance

Using the optimal weights $\hat{\mathbf{w}}(\cdot) = \mathbf{K}^{-1} \mathbf{k}(\mathbf{x})$, the **kriging variance** (MSE of the kriging predictor) **does depend on the process variance**:

$$\begin{aligned}
 \text{MSE}(\hat{Y}_{\mathbf{x}}) &= \mathbb{E} \left[\left(Y_{\mathbf{x}} - \hat{Y}_{\mathbf{x}} \right)^2 \right] \\
 &= \sigma^2 \mathbf{K}(\mathbf{x}, \mathbf{x}) - 2 \hat{\mathbf{w}}(\mathbf{x})^T (\sigma^2 \mathbf{k}(\mathbf{x})) + \hat{\mathbf{w}}(\mathbf{x})^T (\sigma^2 \mathbf{K}) \hat{\mathbf{w}}(\mathbf{x}) \\
 &= \sigma^2 \left(\mathbf{K}(\mathbf{x}, \mathbf{x}) - 2 \mathbf{k}(\mathbf{x})^T \mathbf{K}^{-1} \mathbf{k}(\mathbf{x}) + \mathbf{k}(\mathbf{x})^T \mathbf{K}^{-1} \mathbf{K} \mathbf{K}^{-1} \mathbf{k}(\mathbf{x}) \right) \\
 &= \sigma^2 \left(\mathbf{K}(\mathbf{x}, \mathbf{x}) - \mathbf{k}(\mathbf{x})^T \mathbf{K}^{-1} \mathbf{k}(\mathbf{x}) \right).
 \end{aligned}$$

Kriging Variance

Using the optimal weights $\hat{\mathbf{w}}(\cdot) = \mathbf{K}^{-1} \mathbf{k}(\mathbf{x})$, the **kriging variance** (MSE of the kriging predictor) **does depend on the process variance**:

$$\begin{aligned} \text{MSE}(\hat{Y}_{\mathbf{x}}) &= \mathbb{E} \left[\left(Y_{\mathbf{x}} - \hat{Y}_{\mathbf{x}} \right)^2 \right] \\ &= \sigma^2 \mathbf{K}(\mathbf{x}, \mathbf{x}) - 2 \hat{\mathbf{w}}(\mathbf{x})^T (\sigma^2 \mathbf{k}(\mathbf{x})) + \hat{\mathbf{w}}(\mathbf{x})^T (\sigma^2 \mathbf{K}) \hat{\mathbf{w}}(\mathbf{x}) \\ &= \sigma^2 \left(\mathbf{K}(\mathbf{x}, \mathbf{x}) - 2 \mathbf{k}(\mathbf{x})^T \mathbf{K}^{-1} \mathbf{k}(\mathbf{x}) + \mathbf{k}(\mathbf{x})^T \mathbf{K}^{-1} \mathbf{K} \mathbf{K}^{-1} \mathbf{k}(\mathbf{x}) \right) \\ &= \sigma^2 \left(\mathbf{K}(\mathbf{x}, \mathbf{x}) - \mathbf{k}(\mathbf{x})^T \mathbf{K}^{-1} \mathbf{k}(\mathbf{x}) \right). \end{aligned}$$

As for the MLE criterion, we can use the **optimal profile variance**

$$\sigma_{\text{opt}}^2 = \frac{1}{N} \mathbf{y}^T \mathbf{K}^{-1} \mathbf{y}$$

to see that **the kriging variance is essentially the same as the Golomb–Weinberger criterion**.

Outline

- 1 Deterministic and Statistical Views of Kernel Methods
- 2 Parametrization Criteria
- 3 Computational Aspects**
- 4 Numerical Illustrations

Computing the Cardinal Functions

In the standard basis we find the cardinal basis functions $\ell_j(\mathbf{x}_i) = \delta_{ij}$ via

$$\mathbf{K}\ell(\mathbf{x}) = \mathbf{k}(\mathbf{x}) \iff \ell(\mathbf{x})^T = \mathbf{k}(\mathbf{x})^T \mathbf{K}^{-1},$$

where $\ell(\cdot) = (\ell_1(\cdot) \ \cdots \ \ell_N(\cdot))^T$.

Computing the Cardinal Functions

In the standard basis we find the cardinal basis functions $\ell_j(\mathbf{x}_i) = \delta_{ij}$ via

$$\mathbf{K}\ell(\mathbf{x}) = \mathbf{k}(\mathbf{x}) \iff \ell(\mathbf{x})^T = \mathbf{k}(\mathbf{x})^T \mathbf{K}^{-1},$$

where $\ell(\cdot) = (\ell_1(\cdot) \ \cdots \ \ell_N(\cdot))^T$.

Moreover, in any alternate basis defined by (see [PS11])

$$\mathbf{v}(\mathbf{x})^T = \mathbf{k}(\mathbf{x})^T \mathbf{T}, \quad \mathbf{K} = \mathbf{V}\mathbf{T}^{-1},$$

we have

$$\begin{aligned} \ell(\mathbf{x})^T &= \mathbf{k}(\mathbf{x})^T \mathbf{K}^{-1} \\ &= \mathbf{v}(\mathbf{x})^T \mathbf{T}^{-1} \mathbf{T} \mathbf{V}^{-1} \\ &= \mathbf{v}(\mathbf{x})^T \mathbf{V}^{-1}. \end{aligned}$$

Examples of Alternate Bases

Newton basis [MS09, PS11]: **via Cholesky factorization** $K = NN^T$

$$\mathbf{n}(\mathbf{x})^T = \mathbf{k}(\mathbf{x})^T N^{-T}, \quad V = N, \quad T = N^{-T}$$

SVD basis [PS11, DMS13]: **via SVD** $K = Q\Sigma^2Q^T$

$$\mathbf{v}(\mathbf{x})^T = \mathbf{k}(\mathbf{x})^T Q\Sigma^{-1}, \quad V = Q\Sigma, \quad T = Q\Sigma^{-1}$$

HS-SVD basis [FM12]: **via HS-SVD** $\Psi\Lambda_1\Phi_1^T = K$

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

Remark

Can also use *low rank approximate bases* (see, e.g., [PS11, FM12]).

Computing $\log \det K$

We **compute determinants using logarithms** to prevent underflow errors that may/will arise for small enough values of the shape parameter.

Computing log det K

We **compute determinants using logarithms** to prevent underflow errors that may/will arise for small enough values of the shape parameter.

Mathematically, computing log det K is straightforward and stable using the Hilbert–Schmidt SVD $K = \Psi \Lambda_1 \Phi_1^T$, i.e.,

$$\log \det K = \log \det \Psi + \log \det \Lambda_1 + \log \det \Phi_1^T.$$

Computing log det K

We **compute determinants using logarithms** to prevent underflow errors that may/will arise for small enough values of the shape parameter.

Mathematically, computing log det K is straightforward and stable using the Hilbert–Schmidt SVD $K = \Psi \Lambda_1 \Phi_1^T$, i.e.,

$$\log \det K = \log \det \Psi + \log \det \Lambda_1 + \log \det \Phi_1^T.$$

Computationally,

- the **very small eigenvalues can be handled safely** by taking their logarithms (since Λ_1 is diagonal),
- Φ_1^T gets inverted while forming the stable basis, and
- Ψ gets inverted while computing an interpolant, so the **cost of computing $\log(\det(K))$ is negligible.**



Computing the Power Function

In addition to the ill-conditioning which may be present in the matrix K (and so K^{-1}), **evaluation of the power function is susceptible to numerical cancelation.**

Computing the Power Function

In addition to the ill-conditioning which may be present in the matrix K (and so K^{-1}), **evaluation of the power function is susceptible to numerical cancellation.**

Consider

$$\tilde{K} = \begin{pmatrix} K(\mathbf{x}, \mathbf{x}) & \mathbf{k}(\mathbf{x})^T \\ \mathbf{k}(\mathbf{x}) & K \end{pmatrix},$$

so that

$$\det(\tilde{K}) = \det \left(\begin{pmatrix} 1 & \mathbf{k}(\mathbf{x})^T \\ \mathbf{0}_N & K \end{pmatrix} \begin{pmatrix} K(\mathbf{x}, \mathbf{x}) - \mathbf{k}(\mathbf{x})^T K^{-1} \mathbf{k}(\mathbf{x}) & \mathbf{0}_N^T \\ K^{-1} \mathbf{k}(\mathbf{x}) & I_N \end{pmatrix} \right)$$

Computing the Power Function

In addition to the ill-conditioning which may be present in the matrix K (and so K^{-1}), **evaluation of the power function is susceptible to numerical cancelation.**

Consider

$$\tilde{K} = \begin{pmatrix} K(\mathbf{x}, \mathbf{x}) & \mathbf{k}(\mathbf{x})^T \\ \mathbf{k}(\mathbf{x}) & K \end{pmatrix},$$

so that

$$\begin{aligned} \det(\tilde{K}) &= \det \left(\begin{pmatrix} 1 & \mathbf{k}(\mathbf{x})^T \\ \mathbf{0}_N & K \end{pmatrix} \begin{pmatrix} K(\mathbf{x}, \mathbf{x}) - \mathbf{k}(\mathbf{x})^T K^{-1} \mathbf{k}(\mathbf{x}) & \mathbf{0}_N^T \\ K^{-1} \mathbf{k}(\mathbf{x}) & I_N \end{pmatrix} \right) \\ &= \det(K) (K(\mathbf{x}, \mathbf{x}) - \mathbf{k}(\mathbf{x})^T K^{-1} \mathbf{k}(\mathbf{x})). \end{aligned}$$

Computing the Power Function

In addition to the ill-conditioning which may be present in the matrix K (and so K^{-1}), **evaluation of the power function is susceptible to numerical cancelation.**

Consider

$$\tilde{K} = \begin{pmatrix} K(\mathbf{x}, \mathbf{x}) & \mathbf{k}(\mathbf{x})^T \\ \mathbf{k}(\mathbf{x}) & K \end{pmatrix},$$

so that

$$\begin{aligned} \det(\tilde{K}) &= \det \left(\begin{pmatrix} 1 & \mathbf{k}(\mathbf{x})^T \\ \mathbf{0}_N & K \end{pmatrix} \begin{pmatrix} K(\mathbf{x}, \mathbf{x}) - \mathbf{k}(\mathbf{x})^T K^{-1} \mathbf{k}(\mathbf{x}) & \mathbf{0}_N^T \\ K^{-1} \mathbf{k}(\mathbf{x}) & I_N \end{pmatrix} \right) \\ &= \det(K) (K(\mathbf{x}, \mathbf{x}) - \mathbf{k}(\mathbf{x})^T K^{-1} \mathbf{k}(\mathbf{x})). \end{aligned}$$

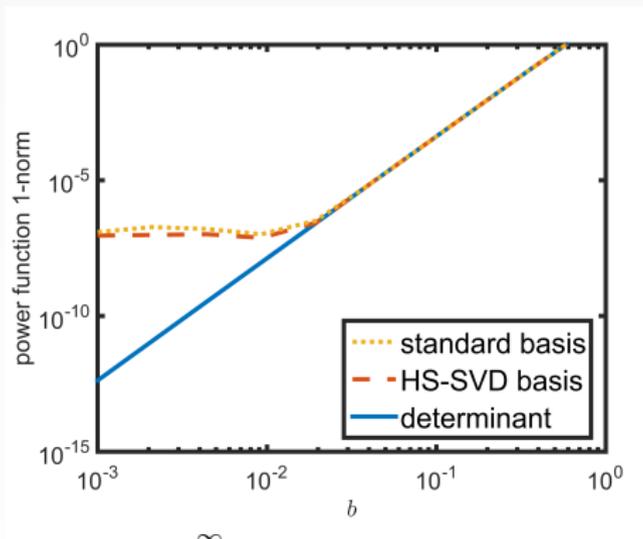
Therefore, the **power function can be computed via** [FWL04, Sch05]

$$P_{K, \mathcal{X}}(\mathbf{x}) = \sqrt{\frac{\det(\tilde{K})}{\det(K)}}.$$

Outline

- 1 Deterministic and Statistical Views of Kernel Methods
- 2 Parametrization Criteria
- 3 Computational Aspects
- 4 Numerical Illustrations**

Computing the Power Function — Example

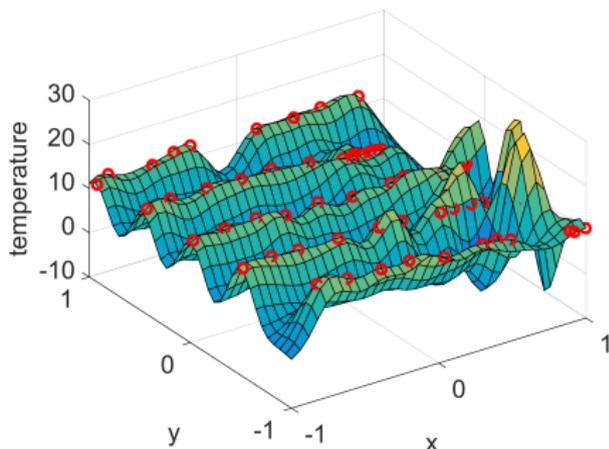
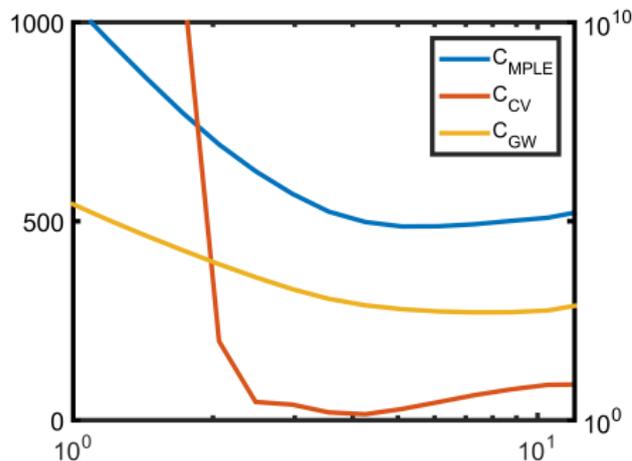


Analytic Chebyshev kernel $K(x, z) = \sum_{n=0}^{\infty} \lambda_n \varphi_n(x) \varphi_n(z)$ on 11 Chebyshev points in $[-1, 1]$

$$\lambda_0 = \frac{1}{2}, \quad \lambda_n = \frac{(1-b)b^n}{2b}, \quad \varphi_n(x) = \sqrt{2 - \delta_{n0}} T_n(x),$$

$$K(x, z) = \frac{1}{2} + (1-b) \frac{b(1-b^2) - 2b(x^2 + z^2) + (1+3b^2)xz}{(1-b^2)^2 + 4b(b(x^2 + z^2) - (1+b^2)xz)}$$

Using various criteria and isotropic kernels to fit track data [Dav14]

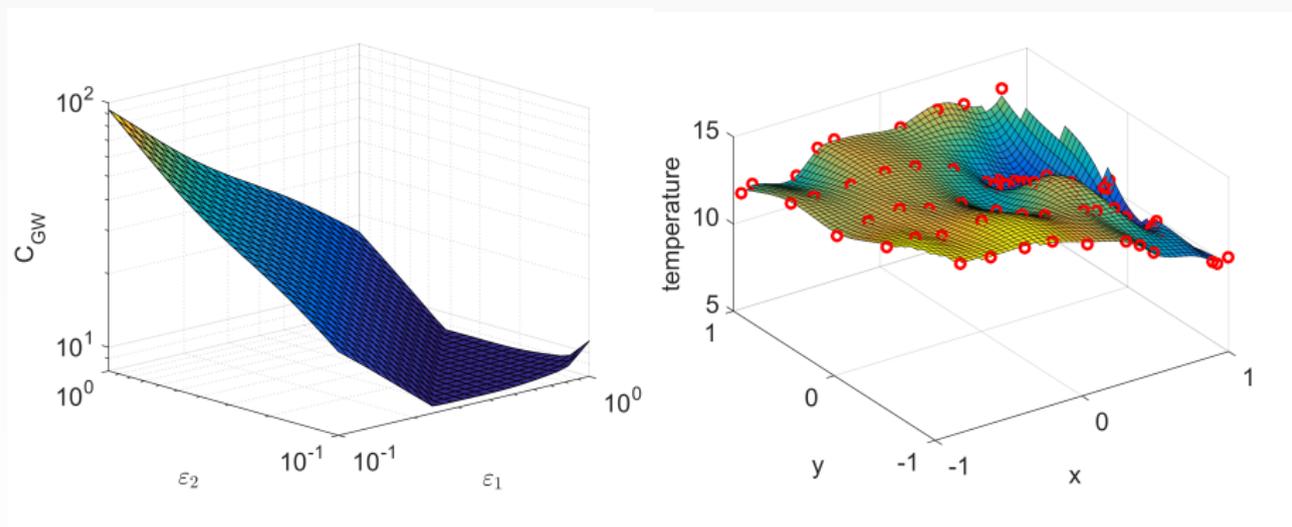


Interpolation with isotropic Gaussian kernel, $\varepsilon = 6$:

$$\varepsilon_{\text{GWopt}} = 7.3162, \quad \varepsilon_{\text{CVopt}} = 4.2518, \quad \varepsilon_{\text{MPLEopt}} = 5.0950$$

Using C_{GW} and anisotropic kernels to fit track data

[Dav14]

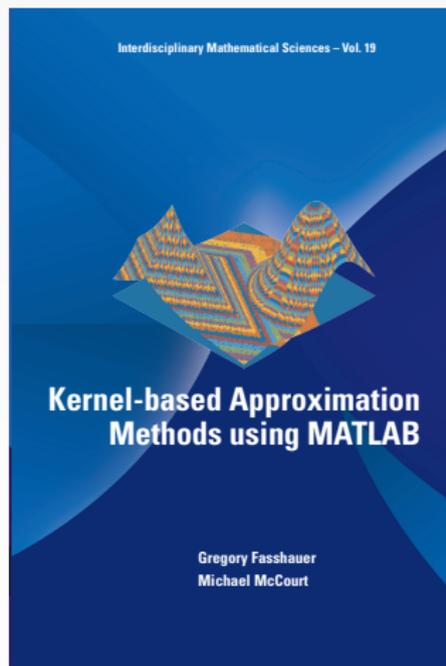


Interpolation with anisotropic tensor product C^2 Matérn kernels:

$$\epsilon_{1opt} = 2.37\epsilon_{2opt}, \quad \epsilon_{1opt} = 0.6803, \quad \epsilon_{2opt} = 0.2875$$

Summary

- Explained various criteria for choosing “optimal” kernel parameters (including C_{GW} , based on error bound)
- Reliable evaluation of these criteria requires
 - alternate (stable) bases
 - kriging variance with process variance
 - determinant formula for power function
- Vast **applications**
 - function interpolation/approximation
 - numerical solution of PDEs (collocation, MFS, MPS, RBF-FD)
 - machine learning (RBF network regression, low-rank approximation for SVM classification)
 - ...



MATLAB code available at

<http://math.iit.edu>

~mccomic/gaussqr



One more inspiring perspective



Calanque

References I

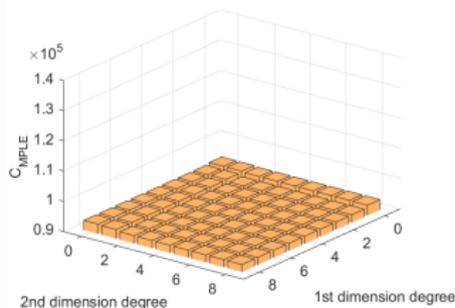
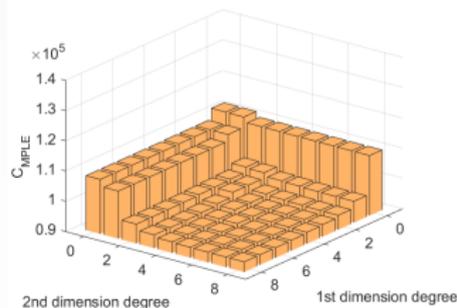
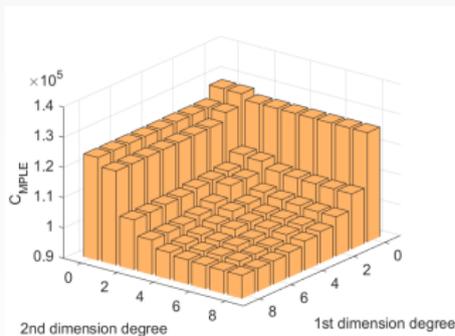
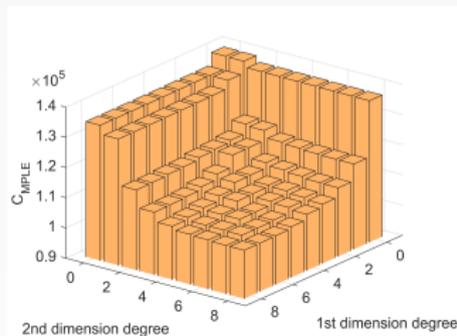
- [Dav14] Oleg Davydov, *Tsfit: A software package for two-stage scattered data fitting*, <https://www.staff.uni-giessen.de/odavydov/tsfit/>, February 2014.
- [DMS13] Stefano De Marchi and Gabriele Santin, *A new stable basis for radial basis function interpolation*, *J. Comput. Appl. Math.* **253** (2013), 1–13.
- [FM12] G. E. Fasshauer and M. J. McCourt, *Stable evaluation of Gaussian radial basis function interpolants*, *SIAM J. Sci. Comput.* **34** (2012), no. 2, A737–A762.
- [FWL04] B. Fornberg, G. Wright, and E. Larsson, *Some observations regarding interpolants in the limit of flat radial basis functions*, *Comput. Math. Appl.* **47** (2004), 37–55.
- [GW59] M. Golomb and H. F. Weinberger, *Optimal approximation and error bounds*, *On Numerical Approximation* (R. E. Langer, ed.), University of Wisconsin Press, 1959, pp. 117–190.
- [MS09] Stefan Müller and Robert Schaback, *A Newton basis for kernel spaces*, *J. Approx. Theory* **161** (2009), no. 2, 645–655.

References II

- [PS11] M. Pazouki and R. Schaback, *Bases for kernel-based spaces*, J. Comput. Appl. Math. **236** (2011), no. 4, 575–588.
- [Sch05] Robert Schaback, *Multivariate interpolation by polynomials and radial basis functions*, Constr. Approx. **21** (2005), 293–317.
- [Sch11] _____, *The missing Wendland functions*, Adv. Comput. Math. **34** (2011), no. 1, 67–81.

Using C_{MPLE} with universal kriging to fit glacier data

[Dav14]



Interpolation with anisotropic Wendland kernels and polynomials:

$$\epsilon = (20 \ 21), \quad \epsilon = (16 \ 14), \quad \epsilon = (8 \ 11), \quad \epsilon = (5 \ 4)$$



Using $C_{MPL E}$ with universal kriging to fit glacier data

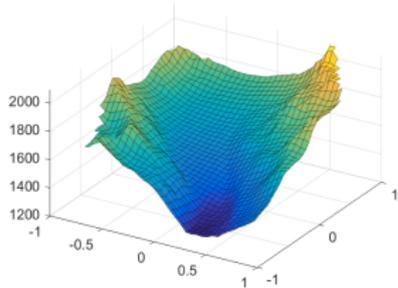
[Dav14]

Interpolation with **anisotropic “missing” Wendland kernels** [Sch11]

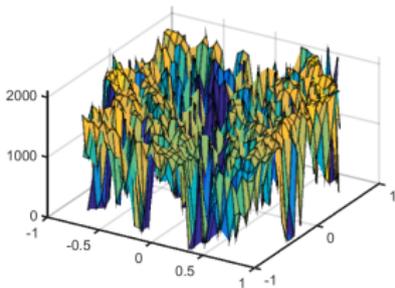
$$\kappa_{3, \frac{3}{2}}(r) \doteq \left(1 - 7r^2 - \frac{81}{4}r^4\right) \sqrt{1 - r^2} - \frac{15}{4}r^4(6 + r^2) \log\left(\frac{r}{1 + \sqrt{1 - r^2}}\right)$$

and **polynomial trend**

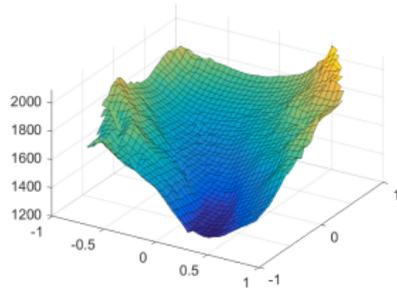
$$s(\mathbf{x}) = \mathbf{k}(\mathbf{x})^T \mathbf{K}^{-1}(\mathbf{y} - \mathbf{P}\beta) + \mathbf{p}(\mathbf{x})^T \beta, \quad \beta = (\mathbf{P}^T \mathbf{K}^{-1} \mathbf{P})^{-1} \mathbf{P}^T \mathbf{K}^{-1} \mathbf{y}.$$



$$\varepsilon = (5 \ 4), \text{deg} = (0 \ 0)$$



$$\varepsilon = (20 \ 21)$$



$$\varepsilon = (20 \ 21), \text{deg} = (8 \ 8)$$

Improved efficiency with hybrid/multiscale methods

ε	degree	density	$K = LL^T$	times (s) $C_{\text{MPLE}}(\varepsilon, \mathbf{p})$	eval
(20 21)	0	0.25%	0.59	0.03	0.32
(20 21)	8		0.59	0.21	0.33
(5 4)	0	4.17%	4.04	0.35	0.63
(5 4)	8		4.04	1.45	0.64

Handle

- large-scale trends with polynomials (or global kernels)
- small-scale features with compactly supported kernels.