



Institut für Numerische Simulation

Rheinische Friedrich-Wilhelms-Universität Bonn

Wegelerstraße 6 • 53115 Bonn • Germany
phone +49 228 73-3427 • fax +49 228 73-7527
www.ins.uni-bonn.de

C. Rieger

**Incremental Kernel Based Approximations for
Bayesian Inverse Problems**

INS Preprint No. 1807

May 2018

INCREMENTAL KERNEL BASED APPROXIMATIONS FOR BAYESIAN INVERSE PROBLEMS

CHRISTIAN RIEGER

ZUSAMMENFASSUNG. We provide an interpretation for the covariance of the predictive process of Bayesian Gaussian process regression as reproducing kernel of a subset of the Cameron Martin space of the prior. We demonstrate that this deterministic viewpoint enables us to relate particular greedy methods using that subset kernel to instances of powerful low-rank matrix approximation techniques such as *adaptive cross approximation* or *pivoted Cholesky decomposition*. In particular, we can show convergence results for such algorithms which appear to be novel in the case of finitely smooth kernels.

Moreover, we consider the inverse problem to reconstruct a parametrized diffusion coefficient from point evaluations of the solution to a diffusion equation with that parametrized coefficient. To this end, we present a Gaussian process regression based approach to approximate the observation operator. The error estimates for this approximation methods are capable to take deterministic model errors explicitly into account. Finally, we show how the findings about incremental low rank approximations can be applied to these reconstruction problems.

1. INTRODUCTION

Gaussian process regression is an important tool to reconstruct a function from finitely many data. The Gaussian process has many different interpretations. We outline connections between Gaussian process regression and numerical approximation methods in reproducing kernel Hilbert spaces (RKHS). Here, we especially focus on the deterministic interpretation of the predictive process and its covariance. It turns out that the covariance of the predictive process is nothing but the so-called power kernel from [10]. The power kernel construction naturally shows up in the *adaptive cross approximation* method, though this powerful method can be formulated in much more generality, using e.g. non-symmetric kernels, see [1]. Algorithmically both methods can be related to the variants of the Cholesky decomposition, see [6]. These findings allow very efficient low-rank approximations of the covariance of the predictive process. In particular, the interpretation in terms of reproducing kernel methods allows to derive an error analysis for the case of finitely smooth covariance functions.

As an application, we apply the general framework of Gaussian process regression to an inverse problem. The inverse problem is to reconstruct a parametrized diffusion coefficient from point evaluations of solutions to the stationary diffusion equation. Here, we follow [15], where a very general theory to bound the reconstruction error for the diffusion coefficients is presented. The difference to the approach in [15] is that we do not assume to have arbitrarily precise solutions to the differential equation. Of course, this is a constructed model problem. If we had all data available

Date: 8. Mai 2018.

to compute the solution to the differential equation, there is no need to recover the diffusion coefficient. So our analysis might be directly applied if the solution to a differential equation is stored as data but the information about the diffusion coefficients are lost. We, however, treat this numerical model error as a prototype for more general types of model errors. Such a model error might be the wrong differential equation to describe the physical processes accurately, or many more. We are convinced that many of those model errors are deterministic in nature and hence, we present an analysis for this case.

The necessity to include the discretization error of the solution to the differential equation is also discussed in a statistical framework in [2]. We pursue a similar approach as we also employ a kernel method to solve the differential equation. As already encountered in Section 2, those kernel methods can be seen as Gaussian process regression. We, however, employ a deterministic a priori error analysis as developed in [4]. To our mind, the deterministic error model is more realistic in the context of model errors. To account for such a deterministic error, we employ techniques from machine learning. Moreover, we can use techniques from low-rank matrix approximation methods to distribute the points according to the covariance of the predictive distribution function in the spirit of [5, 7]. We develop a new error analysis for the inverse reconstruction problem as outlined above. Moreover, we discuss how techniques from adaptive cross approximation can be applied in such problems especially in the design of new observation points.

The remainder of the manuscript is organized as follows: In Section 2, we recall basic facts from the theory of Gaussian process regression and reproducing kernel Hilbert space methods. Moreover, we present a new interpretation of the adaptive cross approximation in the setting of reproducing kernel Hilbert space methods. In Section 3, we present the inverse problem to determine the parametrized diffusion coefficient from point evaluations of the stationary diffusion equation with that coefficient. In Section 4, we present a numerical scheme to solve the differential equation and present its error analysis. Finally, we present an error analysis for the inverse problem in Section 5. We end with some concluding remarks in Section 6.

2. GAUSSIAN PROCESS REGRESSION AND ADAPTIVE CROSS APPROXIMATION

Here, we follow [13, Chapter 2.2]. Let $\mathbb{X} \subset \mathbb{R}^d$. We consider a Gaussian process

$$f \sim \nu_0 = \mathcal{N}(0, \mathcal{C}) \quad \text{or in other notation} \quad f \sim \text{GP}(0, \kappa).$$

Let $\Xi_N = \{\xi_1, \dots, \xi_N\} \subset \mathbb{X}$ be a discrete set. The goal is to extract information about the function f from finitely many point values. We get for the distribution

$$f|_{\Xi_N} \sim \text{GP}(\mathbf{0}, S_{\Xi_N} S_{\Xi_N}^*).$$

In order to formulate the problem more precisely, we need to specify the prior or the Gaussian measure. We consider a symmetric positive definite kernel function $\kappa : \mathbb{X} \times \mathbb{X} \rightarrow \mathbb{R}$. Then, there is a uniquely defined Hilbert space $\mathcal{H}_\kappa(\mathbb{X}) \subset \{f : \mathbb{X} \rightarrow \mathbb{R}\}$ which is called reproducing kernel Hilbert space (RKHS) and satisfies

$$\begin{aligned} \kappa_x &:= \kappa(\cdot, x) \in \mathcal{H}_\kappa(\mathbb{X}) \quad \text{for all } x \in \mathbb{X} \\ f(x) &= (f, \kappa_x)_{\mathcal{H}_\kappa(\mathbb{X})} \quad \text{for all } f \in \mathcal{H}_\kappa(\mathbb{X}) \text{ and for all } x \in \mathbb{X}. \end{aligned}$$

The integral operator

$$(1) \quad \mathcal{C}_\kappa : L_2(\mathbb{X}) \rightarrow L_2(\mathbb{X}), \quad f \mapsto \mathcal{C}_\kappa(f) = \int_{\mathbb{X}} \kappa(\cdot, y) f(y) dy$$

gives rise to a (centered) Gaussian measure $\nu_0 = \mathcal{N}(0, \mathcal{C}_\kappa)$ on $L_2(\mathbb{X})$. Furthermore we have

$$\text{Range}(\mathcal{C}_\kappa^{\frac{1}{2}}) = \mathcal{H}_\kappa(\mathbb{X}) \subset L_2(\mathbb{X}),$$

i.e., the RKHS is the Cameron Martin space on ν_0 . This means that we choose ν_0 as prior. We consider the case that $\nu_0(C(\mathbb{X})) = 1$ which requires some regularity of the kernel $\kappa : \mathbb{X} \times \mathbb{X} \rightarrow \mathbb{R}$ and is in our setting usually related to Sobolev embeddings. We introduce the maps

$$\begin{aligned} S_{\Xi_N} : \mathcal{H}_\kappa(\mathbb{X}) &\rightarrow \mathbb{R}^N, & f &\mapsto (f(\xi_1), \dots, f(\xi_N))^\top \\ S_{\Xi_N}^* : \mathbb{R}^N &\rightarrow \mathcal{H}_\kappa(\mathbb{X}), & \mathbf{c} &\mapsto \sum_{\xi_j \in \Xi_N} c_j \kappa_{\xi_j}. \end{aligned}$$

The notation is justified since we have

$$(\mathbf{d}, S_{\Xi_N} f)_{\ell_2(\Xi_N)} = \sum_{j=1}^N d_j f(\xi_j) = (f, S_{\Xi_N}^* \mathbf{d})_{\mathcal{H}_\kappa(\mathbb{X})} \quad \text{for all } f \in \mathcal{H}_\kappa(\mathbb{X}) \text{ and all } \mathbf{d} \in \mathbb{R}^N.$$

We use the notation

$$\mathcal{I}_{\Xi_N} : \mathcal{H}_\kappa(\mathbb{X}) \rightarrow \mathcal{H}_\kappa(\mathbb{X}), \quad f \mapsto \mathcal{I}_{\Xi_N;f}(\cdot) = \left(S_{\Xi_N}^* (S_{\Xi_N} S_{\Xi_N}^*)^{-1} S_{\Xi_N}(f) \right) (\cdot).$$

Moreover, we observe that

$$S_{\Xi_N} S_{\Xi_N}^* = \mathbf{K}_{\Xi_N, \Xi_N} = \begin{pmatrix} \kappa(\xi_1, \xi_1) & \dots & \kappa(\xi_1, \xi_N) \\ \vdots & \ddots & \vdots \\ \kappa(\xi_N, \xi_1) & \dots & \kappa(\xi_N, \xi_N) \end{pmatrix} \in \mathbb{R}^{N \times N}$$

is a symmetric and positive definite matrix. We introduce the space

$$\ker(S_{\Xi_N}) = \{f \in \mathcal{H}_\kappa(\mathbb{X}) : S_{\Xi_N}(f) = \mathbf{0}\} \subset \mathcal{H}_\kappa(\mathbb{X})$$

and note that its reproducing kernel is given by [10, Section 5]

$$(2) \quad \kappa_{\Xi_N} : \mathbb{X} \times \mathbb{X} \rightarrow \mathbb{R}, \quad (x, y) \mapsto \left((\text{Id} - \mathcal{I}_{\Xi_N})^{(1)} (\text{Id} - \mathcal{I}_{\Xi_N})^{(2)} \kappa(\cdot, \cdot) \right) (x, y),$$

where the superscript denotes to which argument the function is applied. Following [10], we call the kernel (2) the *power kernel*. The terminology is motivated by the identity

$$(3) \quad \kappa_{\Xi_N}(x, x) = P_{\Xi_N}^2(x) = \text{dist}_{\mathcal{H}_\kappa^*(\mathbb{X})}^2(\delta_x, \{\delta_\xi : \xi \in \Xi_N\}),$$

where P_{Ξ_N} is usually called the *power function*. Then, we have the following representation

$$f(\cdot) | S_{\Xi_N}(f) \sim \text{GP}(\mathcal{I}_{\Xi_N;f}(\cdot), \kappa_{\Xi_N}(\cdot, \cdot))$$

for the predictive process, see [13, Eq. (2.19)] or [15, Eq. (3.6)].

2.1. Recursive interpolation. We assume that the data points are not given as a single set of points but that they can be viewed as a data stream. We denote by

$$\Xi_N = \{\xi_1, \dots, \xi_{N(\ell)}\} = \Xi_{N-1} \cup \{\xi_N\} \quad \text{and} \quad \Xi_0 = \emptyset.$$

The first observation is that such a data stream fits the asymptotic nature of most error estimates better than a single set of points. The perspective of incremental approximations can also be employed in a low-rank approximation of the kernel matrix. Following [10, Section 5], we have the following recursion for the power kernel (2)

$$\begin{aligned} \kappa_{\Xi_0, \Xi_0}(x, y) &= \kappa(x, y) \quad \text{and} \\ (4) \quad \kappa_{\Xi_{N+1}, \Xi_{N+1}}(x, y) &= \kappa_{\Xi_N, \Xi_N}(x, y) - \frac{\kappa_{\Xi_N, \Xi_N}(x, \xi_{N+1})\kappa_{\Xi_N, \Xi_N}(\xi_{N+1}, y)}{\kappa_{\Xi_N, \Xi_N}(\xi_{N+1}, \xi_{N+1})} \quad \text{for all } x, y \in \mathbb{X}. \end{aligned}$$

Such a recursion appears also for the symmetric kernel case in the adaptive cross approximation, c.f. [1, Eq 16]. Following [10, Section 5], we consider also the matrix version of that recursion. We define the kernel matrix for the n -th power kernel defined in (4) with $n \leq N$

$$\mathbb{K}_{n; \Xi_n, \Xi_n} = \begin{pmatrix} \kappa_{\Xi_n, \Xi_n}(\xi_1, \xi_1) & \dots & \kappa_{\Xi_n, \Xi_n}(\xi_1, \xi_N) \\ \vdots & \ddots & \vdots \\ \kappa_{\Xi_n, \Xi_n}(\xi_N, \xi_1) & \dots & \kappa_{\Xi_n, \Xi_n}(\xi_N, \xi_N) \end{pmatrix} \in \mathbb{R}^{N \times N},$$

evaluated at the point set Ξ_n . This is the matrix C_j from [10, Eq. (5)]. As outlined there, this matrix is the n -th iterate of the *inline Cholesky decomposition*

$$\begin{aligned} \mathbb{K}_{0; \Xi_n, \Xi_n} &= \mathbf{K}_{\Xi_n, \Xi_n} \quad \text{and} \\ (5) \quad \mathbb{K}_{n; \Xi_n, \Xi_n} &= \mathbb{K}_{n-1; \Xi_n, \Xi_n} - \frac{1}{\mathbb{K}_{n-1; \Xi_n, \Xi_n}(n, n)} \mathbb{K}_{n-1; \Xi_n, \Xi_n}(:, n) \otimes \mathbb{K}_{n-1; \Xi_n, \Xi_n}(n, :), \end{aligned}$$

using the same MATLAB inspired notation as in [10, Eq. (6)]. Algorithmically, this relates to the recursion of the adaptive cross approximation as in [1, Section 3.1] for symmetric kernels. The iteration (5) yields a Cholesky decomposition of the kernel matrix, i.e.,

$$\mathbf{K}_{\Xi_n, \Xi_n} = \mathbb{K}_{N; \Xi_n, \Xi_n} \mathbb{K}_{N; \Xi_n, \Xi_n}^\top.$$

As outlined in [12, Section 6], there is a relation between matrix factorizations and special bases, for instance a Cholesky decomposition of the kernel matrix leads to a $(\cdot, \cdot)_{\mathcal{H}_\kappa(\mathbb{X})}$ -orthogonal basis. In order to describe this, we define the vector-field

$$\boldsymbol{\kappa}_{\Xi_n} := S_{\Xi_n}^{(2)} \kappa(\cdot, \cdot) = (\kappa(\cdot, \xi_1), \dots, \kappa(\cdot, \xi_N))^\top$$

and construct a new *Netwon-type* basis as

$$(6) \quad (\psi_{1; \Xi_n}, \dots, \psi_{N; \Xi_n}) = S_{\Xi_n}^{(2)} \kappa(\cdot, \cdot) \cdot \left(\mathbb{K}_{N; \Xi_n, \Xi_n}^\top \right)^{-1},$$

where $\mathbb{K}_{N; \Xi_n, \Xi_n}^\top$ is by construction via the Cholesky decomposition an upper triangular matrix. We have that

$$\text{span} \{ \psi_{1; \Xi_n}, \dots, \psi_{N; \Xi_n} \} = \text{span} \{ \kappa_{\xi_1}, \dots, \kappa_{\xi_N} \} =: V_{\Xi_n}.$$

The basis ψ_j was already discussed in [11] with a different normalization. The main feature of this basis is that 'old' basis elements remain unchanged if a new point is added. This makes this new basis especially favorable in sequential applications. A related basis change is also described in [1, Lemma 4], with the choices $\hat{\ell}_i = \kappa_{\xi_i}$ and $C = \mathbb{K}_{N; \Xi_N, \Xi_N}^\top$.

2.2. Low-rank approximation of kernel matrices. We consider a very fine discrete set of points $\Upsilon_{N_\infty} \subset \mathbb{X}$ with $N_\infty < \infty$. In theory, it would be possible to build the matrix $\mathbf{K}_{\Upsilon_{N_\infty}, \Upsilon_{N_\infty}} \in \mathbb{R}^{N_\infty \times N_\infty}$ but N_∞ is assumed to be so large that direct inversion of this matrix should be avoided. Low-rank approximation can provide reliable approximations to that matrix which are cheaper to store and allow for faster matrix algorithms. To this end, we fix a precision $\varepsilon_{\text{chol}} > 0$ and try to find a rank $N_{\varepsilon_{\text{chol}}}$ and a matrix $\mathbf{A} \in \mathbb{R}^{N_\infty \times N_\infty}$ such that

$$\|\mathbf{K}_{\Upsilon_{N_\infty}, \Upsilon_{N_\infty}} - \mathbf{A}\| \leq \varepsilon_{\text{chol}}$$

and $\text{Rank}(\mathbf{A}) \leq N_{\varepsilon_{\text{chol}}}$. We can use [12, Eq. (20)] to get

$$(7) \quad \kappa(x, y) = \lim_{h_{\Xi_N, \mathbb{X}} \rightarrow 0} \sum_{j=1}^N \psi_{j; \Xi_N}(x) \psi_{j; \Xi_N}(y),$$

where the limit means that the point set Ξ_N needs to get dense in \mathbb{X} which in turn necessarily implies $N \rightarrow \infty$. Moreover, we have the identity (c.f. [12, p. 584])

$$\kappa(x, y) - \kappa_{\Xi_N, \Xi_N}(x, y) = \sum_{j=1}^N \psi_{j; \Xi_N}(x) \psi_{j; \Xi_N}(y).$$

The representation of the kernel (7) leads to

$$(8) \quad \mathbf{K}_{\Upsilon_{N_\infty}, \Upsilon_{N_\infty}} = \lim_{h_{\Xi_N, \mathbb{X}} \rightarrow 0} \sum_{j=1}^N S_{\Upsilon_{N_\infty}} \psi_{j; \Xi_N} \otimes S_{\Upsilon_{N_\infty}} \psi_{j; \Xi_N}.$$

Such a series expansion motivates a choice for the matrix \mathbf{A} as

$$(9) \quad \mathbf{A} = \mathbf{A}_{\Xi_N} = \sum_{j=1}^N S_{\Upsilon_{N_\infty}} \psi_{j; \Xi_N} \otimes S_{\Upsilon_{N_\infty}} \psi_{j; \Xi_N}.$$

The important point which remains open is how to construct the point sets Ξ_n iteratively to achieve a small $\varepsilon_{\text{chol}}$. Both [1] and [6] provide efficient algorithms to compute such a low rank approximation. As already mentioned in [6, Section 3], both approaches are similar if a *total pivoting* strategy is used in [1]. We will present an other kernel-based derivation for such algorithms. In particular, the kernel based approach also allows for an error analysis in the case of finitely smooth kernels.

Kernel-based power greedy method. For any discrete set $\Xi \subset \mathbb{X}$, there is an error formula for kernel-based interpolation

$$|(f - \mathcal{I}_\Xi(f))(x)| \leq P_\Xi(x) \|f\|_{\mathcal{H}_\kappa(\mathbb{X})} \quad \text{for all } f \in \mathcal{H}_\kappa(\mathbb{X}),$$

using the power function from (3), see [16]. The important feature of this error estimate is that the power function does not depend on the specific function f . Hence, such a error estimate is universal for the whole function space $\mathcal{H}_\kappa(\mathbb{X})$. By definition, we observe

$$\kappa_{\Xi}(\xi, \xi) = P_{\Xi}^2(\xi) = 0 \quad \text{for all } \xi \in \Xi.$$

The *power greedy methods* is now defined as follows: Start with a point $\xi_1 \in \mathbb{X}$ and set $\Xi_1 = \{\xi_1\}$. Define iteratively

$$(10) \quad \xi_{n+1} := \arg \max_{x \in \mathbb{X}} P_{\Xi_n}^2(x) \quad \text{and set} \quad \Xi_{n+1} = \Xi_n \cup \{\xi_{n+1}\}.$$

The rationale behind this algorithm is to try to reduce the high values of the power function as quickly as possible and hence to make the power function globally small. For the convergence analysis we refer to [14] for the case of finitely smooth kernels and to [6] for infinitely smooth kernels. As finitely smooth reproducing kernels, we will consider *radial basis functions* with a specific decay of their Fourier transforms. Precisely, we consider translation-invariant kernels Φ such that there is a ϕ satisfying

$$\Phi : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}, x \mapsto \Phi(x - y) = \phi(\|x - y\|_2),$$

where $\phi : \mathbb{R} \rightarrow \mathbb{R}$. The decay condition on the Fourier transform reads: there are constants $0 < c_\phi, C_\phi < \infty$ such that for $\lfloor \tau \rfloor > \frac{d}{2}$ it holds that

$$(11) \quad c_\phi \left(1 + \|\omega\|_2^2\right)^\tau \leq \hat{\Phi}(\omega) \leq C_\phi \left(1 + \|\omega\|_2^2\right)^\tau \quad \text{for all } \omega \in \mathbb{R}^d.$$

We assume also $\phi \in C^2(\mathbb{R})$. This is for instance satisfied for certain Wendland kernels, see [16, proof of Theorem 11.17]. For those kernels, it turns out that

$$\mathcal{H}_\Phi(\mathbb{R}^d) \cong W_2^\tau(\mathbb{R}^d).$$

Moreover, we obtain [16, Corollary 10.48] that

$$\mathcal{H}_\Phi(\mathbb{X}) \cong W_2^\tau(\mathbb{X})$$

if $\mathbb{X} \subset \mathbb{R}^d$ has a Lipschitz boundary and we use a suitable norm which is equivalent to the standard norm on the Sobolev space. Furthermore, the covariance operator (1) behaves in the assumed way [16, Lemma 10.27].

Convergence of P-greedy algorithm for finitely smooth kernels. The algorithm (10) is numerically infeasible since the maximum over the whole set \mathbb{X} is not computable. Following [14, Section 4], we replace the set \mathbb{X} by the finite set $\Upsilon_{N_\infty} = \{v_1, \dots, v_{N_\infty}\}$. The algorithm then reads

$$(12) \quad \begin{aligned} \xi_1 &:= \arg \max_{x \in \Upsilon_{N_\infty}} \kappa(x, x) \quad \text{and} \quad \Xi_1 = \{\xi_1\} \\ \xi_{n+1} &:= \arg \max_{x \in \Upsilon_{N_\infty} \setminus \Xi_n} P_{\Xi_n}^2(x) \quad \text{and} \quad \Xi_{n+1} = \Xi_n \cup \{\xi_{n+1}\}. \end{aligned}$$

Then, we have by [14, Theorem 4.1]

$$(13) \quad \|P_{\Xi_n}\|_{\ell_\infty(\Upsilon_{N_\infty})} \leq C_P n^{-\frac{\tau}{d} + \frac{1}{2}}.$$

This result is enough to get an error estimate for the low-rank matrix approximation, where the error is measured in the Frobenius norm.

Theorem 2.1. *Let $\Xi \subset \mathbb{R}^d$ satisfy an interior cone condition and let $\Phi : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ be of the form (11). Then, there is a constant $c_\epsilon > 0$ such that for all sets Ξ_n created by the discrete power greedy method (12), the matrix \mathbf{A}_{Ξ_n} satisfies the estimate*

$$\|\mathbf{K}_{\Upsilon_{N_\infty}, \Upsilon_{N_\infty}} - \mathbf{A}_{\Xi_n}\|_F \leq C_P N_\infty n^{-\frac{\tau}{d} + \frac{1}{2}}.$$

Proof. We observe using (8) and (9) that

$$(14) \quad \mathbf{K}_{\Upsilon_{N_\infty}, \Upsilon_{N_\infty}} - \mathbf{A}_{\Xi_n} = \begin{pmatrix} \kappa_{\Xi_n, \Xi_n}(v_1, v_1) & \cdots & \kappa_{\Xi_n, \Xi_n}(v_1, v_{N_\infty}) \\ \vdots & \ddots & \vdots \\ \kappa_{\Xi_n, \Xi_n}(v_{N_\infty}, v_1) & \cdots & \kappa_{\Xi_n, \Xi_n}(v_{N_\infty}, v_{N_\infty}) \end{pmatrix} \in \mathbb{R}^{N_\infty \times N_\infty}.$$

This error representation (14) implies an error bound

$$\|\mathbf{K}_{\Upsilon_{N_\infty}, \Upsilon_{N_\infty}} - \mathbf{A}_{\Xi_n}\|_F = \left(\sum_{i=1}^{N_\infty} \sum_{j=1}^{N_\infty} \kappa_{\Xi_n, \Xi_n}^2(v_i, v_j) \right)^{\frac{1}{2}} \leq N_\infty \|P_{\Xi_n}\|_{\ell_\infty(\Upsilon_{N_\infty})}.$$

Now, (13) yields the claim. \square

In order to derive kernel-based error estimates, we also have to study the distribution of the point sets Ξ_n . In particular, we have to bound $h_{\Xi_n, \mathbb{X}}$ in terms of n . To this end, we will provide a new proof for the corollary [14, Corollary 4.3] in order to make the dependence on the finite set Υ_{N_∞} more explicitly. In particular, we need to make sure that the finite set of points Υ_{N_∞} is fine enough to yield

$$\|P_{\Xi_n}\|_{\ell_\infty(\Upsilon_{N_\infty})} \leq \|P_{\Xi_n}\|_{L_\infty(\mathbb{X})} \leq C \|P_{\Xi_n}\|_{\ell_\infty(\Upsilon_{N_\infty})}$$

for a suitable constant $C > 0$.

Theorem 2.2. *Let $\mathbb{X} \subset \mathbb{R}^d$ satisfy an interior cone condition and let $\Phi : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ be of the form (11). Let $h_{\Upsilon_{N_\infty}, \mathbb{X}} \leq h_0$, where h_0 will be defined in (17). Then, for every $\epsilon \in (0, 1)$ there is a constant $c_\epsilon > 0$ such that for all sets Ξ_n created by the discrete power greedy method (12), we have that there is a constant $\theta = \tau - \frac{d}{2} > 0$ such that*

$$(15) \quad h_{\Xi_n, \mathbb{X}} := \sup_{x \in \mathbb{X}} \min_{\xi \in \Xi_n} \|x - \xi\|_2 \leq c_\epsilon n^{-\frac{\theta}{\theta + \epsilon} \frac{1}{d}} \quad \text{for all } .$$

Proof. We mainly follow [14, Corollary 4.3]. Here, we have to use an approximation argument in order to apply [9, Theorem 3.1]. Precisely, [9, Theorem 3.1] provides for every $\epsilon > 0$ the existence of a constant $C(\epsilon)$ such that

$$h_{\Xi_n, \mathbb{X}} \leq C(\epsilon) \|P_{\Xi_n}\|_{L_\infty(\mathbb{X})}^{\frac{1}{\tau - \frac{d}{2} + \epsilon}}.$$

In order to use (13), we have to estimate $\|P_{\Xi_n}\|_{L_\infty(\mathbb{X})}$ in terms of $\|P_{\Xi_n}\|_{\ell_\infty(\Upsilon_{N_\infty})}$. From the condition (11), we have $\lceil \tau \rceil > \frac{d}{2}$. Here, we use [16, Theorem 2.6], to get

$$(16) \quad \|P_{\Xi_n}\|_{L_\infty(\mathbb{X})} \leq C_S \left(h_{\Upsilon_{N_\infty}, \mathbb{X}}^{1-\frac{d}{2d}} \|P_{\Xi_n}\|_{W_{2d}^1(\mathbb{X})} + \|P_{\Xi_n}\|_{\ell_\infty(\Upsilon_{N_\infty})} \right).$$

We use [9, Lemma 4.2] to get

$$|\partial_k P_{\Xi_n}(x)| \leq 2P_{\Xi_n}(x) \sqrt{\partial_k^{(1)} \partial_k^{(2)} \kappa(x, x)} \quad \text{for all } x \in \mathbb{X} \text{ and } 1 \leq k \leq d,$$

where the integrability follows from the assumed regularity of the kernel. Hence, we obtain

$$\|P_{\Xi_n}\|_{W_{2d}^1(\mathbb{X})} \leq C(\kappa, d, \mathbb{X}) \|P_{\Xi_n}\|_{L_\infty(\mathbb{X})}.$$

This implies that if $h_{\Upsilon_{N_\infty}, \mathbb{X}} \leq h_0$ with

$$(17) \quad h_0 \leq (2C_S C(\kappa, d, \mathbb{X}))^{-2}$$

we have by (16)

$$\|P_{\Xi_n}\|_{L_\infty(\mathbb{X})} \leq 2C_S \|P_{\Xi_n}\|_{\ell_\infty(\Upsilon_{N_\infty})} \leq 2C_S C_P n^{-\frac{\tau}{d} + \frac{1}{2}}$$

where the last inequality is given by (13). Hence, we obtain

$$h_{\Xi_n, \mathbb{X}} \leq C(\epsilon) \|P_{\Xi_n}\|_{L_\infty(\mathbb{X})}^{\frac{1}{\tau - \frac{d}{2} + \epsilon}} \leq C(\epsilon) \left(2C_S C_P n^{-\frac{\tau}{d} + \frac{1}{2}} \right)^{\frac{1}{\tau - \frac{d}{2} + \epsilon}}.$$

Manipulating the exponent yields

$$\left(-\frac{\tau}{d} + \frac{1}{2} \right) \frac{1}{\tau - \frac{d}{2} + \epsilon} = -\frac{1}{d} \frac{\tau - \frac{d}{2}}{\tau - \frac{d}{2} + \epsilon}.$$

Setting now $\theta = \tau - \frac{d}{2}$ finishes the proof. \square

The proof also shows the following corollary which is a stronger version of (13).

Corollary 2.3. *Let $\mathbb{X} \subset \mathbb{R}^d$ satisfy an interior cone condition and let $\Phi : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ be of the form (11). Let $h_{\Upsilon_{N_\infty}, \mathbb{X}} \leq h_0$, where h_0 is defined in (17). Then, for ever $\epsilon \in (0, 1)$ there is a constant $c_\epsilon > 0$ such that for all sets Ξ_n created by the discrete power greedy method (12), we have that*

$$\|P_{\Xi_n}\|_{L_\infty(\mathbb{X})} \leq c_\epsilon n^{-\frac{\tau}{d} + \frac{1}{2}}$$

Summarizing, we get results on the distribution of the set of points created by the power greedy method.

3. INVERSE MODEL PROBLEM

Let $D \subset \mathbb{R}^d$ be an open bounded set with boundary ∂D . Let $\tau \in \mathbb{N}$ with $\tau > \frac{d}{2} + 2$. We will assume $\partial D \in C^{\tau+2}$. Let $\phi_k \in C^{\tau+1}(\bar{D})$ for $1 \leq k \leq K$ be given functions with $\|\phi_k\|_{L^\infty(D)} = 1$ for all $1 \leq k \leq K$. Let

$$\mathcal{A} : I^K := [-1, 1]^K \rightarrow C^{\tau+1}(\bar{D}), \quad \mathbf{c} = (c_1, \dots, c_d) \mapsto \mathcal{A}(\mathbf{c}, x) = 1 + \sum_{k=1}^K \frac{c_k}{2K} \phi_k.$$

We directly observe that

$$(18) \quad A_{\min} = \frac{1}{2} \leq \mathcal{A}(\mathbf{c}, x) \leq A_{\max} = \frac{3}{2} \quad \text{for all } \mathbf{c} \in I^K \text{ and all } x \in D.$$

We consider the second order elliptic boundary value problem for given $g \in W_2^{\tau-2}(D)$ with

$$(19) \quad \begin{cases} L_{\mathbf{c}}(u) := -\nabla \cdot (\mathcal{A}(\mathbf{c}, x) \mathbf{Id}_{d \times d} \nabla u(\mathbf{c}, x)) = g(x) & \text{for all } x \in D \\ u(\mathbf{c}, x) = 0 & \text{for all } x \in \partial D \end{cases}$$

and for fixed parameter $\mathbf{c} \in I^K$. We introduce the operator

$$(20) \quad \mathcal{L}_{D;\mathbf{c}} : W_2^\tau(D) \rightarrow W_2^{\tau-2}(D) \times W_2^{\tau-\frac{1}{2}}(\partial D), \quad u \mapsto (L_{\mathbf{c}}u, \gamma_0(u)),$$

where γ_0 is the usual trace map. Elliptic regularity theory [3, Sec. 6.3, Thm. 5] implies that $\mathcal{L}_{D;\mathbf{c}}$ given as in (20) is continuously invertible, i.e., there is a map

$$(21) \quad \mathcal{L}_{D;\mathbf{c}}^{-1} : W_2^{\tau-2}(D) \times W_2^{\tau-\frac{1}{2}}(\partial D) \rightarrow W_2^\tau(D), \quad (g, 0) \mapsto u = \mathcal{L}_{D;\mathbf{c}}^{-1}(g, 0),$$

and we have

$$\|u\|_{W_2^\tau(D)} = \left\| \mathcal{L}_{D;\mathbf{c}}^{-1}(g, 0) \right\|_{W_2^\tau(D)} \leq C(A_{\min}, A_{\max}) \|g\|_{W_2^{\tau-2}(D)},$$

where we used to notation $C(A_{\min}, A_{\max})$ to denote that the constant depends on \mathbf{c} only via the upper and lower bounds in (18). The assumption on the smoothness of D also ensures the existence of a bounded extension operator

$$\mathcal{E}_D : W_2^\tau(D) \rightarrow W_2^\tau(\mathbb{R}^d)$$

such that

$$\mathcal{E}_D u|_D = u \quad \text{and} \quad \|\mathcal{E}_D u\|_{W_2^\tau(\mathbb{R}^d)} \leq \|\mathcal{E}_D\| \|u\|_{W_2^\tau(D)} \quad \text{for all } u \in W_2^\tau(D).$$

Following [15], we consider the mapping

$$G : I^K \rightarrow W_2^\tau(D), \quad \mathbf{c} \mapsto u = \mathcal{L}_{D;\mathbf{c}}^{-1}(g, 0),$$

where $u \in W_2^\tau(D)$ is a solution to (19). Furthermore, we introduce a sampling operator associated with a discrete set $X_N := \{x_1, \dots, x_N\} \subset D \subset \mathbb{R}^d$ as

$$S_{X_N} : W_2^\tau(\mathbb{R}^d) \rightarrow \mathbb{R}^N,$$

which is well-defined due to the condition $\tau > \frac{d}{2}$. The concatenation of the operators

$$(22) \quad \mathcal{G} := S_{X_N} \circ \mathcal{E}_D \circ G : I^K \rightarrow \mathbb{R}^N, \quad \mathbf{c} \mapsto S_{X_N}(\mathcal{E}_D(u)) = S_{X_N}(\mathcal{E}_D \circ \mathcal{L}_{D;\mathbf{c}}^{-1}(g, 0))$$

is called the observation operator in [15]. As in [15], we consider the associated inverse problem to find $\mathbf{c} \in I^K$ from data

$$(23) \quad \mathbf{y}_m = \mathcal{G}(\mathbf{c}_m) + \boldsymbol{\epsilon}_{m;L} \quad \text{for } 1 \leq m \leq M,$$

where $\boldsymbol{\epsilon}_{m;L}$ denotes the numerical model error. Since we cannot evaluate the observation operator for all parameters $\mathbf{c} \in I^K$. We define the set $C_M := \{c_1, \dots, c_M\} \subset I^K$. Hence, we have only the finitely many values

$$(24) \quad S_{C_M}(\hat{\mathcal{G}}) = \left(\hat{\mathcal{G}}(c_1), \dots, \hat{\mathcal{G}}(c_M) \right)^\top \in \mathbb{R}^{N \times M},$$

where $\hat{\mathcal{G}} \approx \mathcal{G}$ is a suitable approximation to the observation operator, at hand to gain information on the continuous observation operator $\mathcal{G} : I^K \rightarrow \mathbb{R}^N$. The strategy of [15] is to assume a prior measure μ_0 such that \mathbf{c} is distributed according to μ_0 . The information, we would like to infer about the parameters is the posterior distribution $\mu^{\mathbf{y}}$ which is the distribution of the conditioned random variable $\mathbf{c}|\mathbf{y}$. Bayes' theorem in this situation reads as follows, see [15, Proposition 2.1].

Proposition 3.1. *Given that $\mathcal{G} : I^K \rightarrow \mathbb{R}^N$ is continuous and we have $\mu_0(I^K) = 1$. Then the distribution of the conditioned random variable $\mathbf{c}|\mathbf{y}$ is absolutely continuous with respect to the prior μ_0 and hence the Radon-Nikodým derivative exists and is given as*

$$(25) \quad \frac{d\mu^{\mathbf{y}}}{d\mu_0}(\mathbf{c}) = \frac{1}{\mathbb{E}_{\mu_0} [\exp(-\|\mathbf{y} - \mathcal{G}(\mathbf{c})\|_2)]} \exp(-\|\mathbf{y} - \mathcal{G}(\mathbf{c})\|_2)$$

Following [15], we construct an approximation to the measure (25) by approximation to the observation operator (22). We construct an approximation

$$(26) \quad \mathcal{G}_{N,M} \approx \mathcal{G}$$

and use this to construct the posterior measure as

$$(27) \quad \frac{d\mu^{\mathbf{y}}}{d\mu_0}(\mathbf{c}) \approx \frac{1}{\mathbb{E}_{\mu_0} [\exp(-\|\mathbf{y} - \mathcal{G}(\mathbf{c})\|_2)]} \exp(-\|\mathbf{y} - \mathcal{G}_{N,M}(\mathbf{c})\|_2)$$

The problem is that $\mathcal{G}(\mathbf{c})$ involves two approximation processes. The first approximation stems from the use of only finitely many data C_M and the second approximation process involves the solution of the system of equations (19) for a fixed parameter.

4. APPROXIMATION OF THE FORWARD PROBLEM – NUMERICAL MODEL ERROR

The forward problem consists in solving the system of equations (19) for a given $\mathbf{c} \in I^K$. As outlined in [2] there is an uncertainty on the solution $u \in W_2^1(D)$ due to the fact that most numerical schemes use only a finite number of point-values of g for the reconstruction of $u = \mathcal{L}_{D;\mathbf{c}}^{-1}(g, 0)$. We treat this uncertainty, in contrast to [2], as a deterministic (numerical) model error. We note as already mentioned in the introduction that this numerical model error serves as prototype for more general model errors.

Following [4], we apply a symmetric kernel-based collocation approach. This collocation method is a special case of generalized kernel-based interpolation. In an abstract framework the generalized kernel-based interpolation can be described as follows. Let $\mathcal{H}_\kappa(\mathbb{X})$ be an RKHS with kernel κ . Consider a set of linear functionals

$$\Lambda_N := \{\lambda_1, \dots, \lambda_N\} \subset \mathcal{H}_\kappa^*(\mathbb{X}),$$

where we assume that $\dim \text{span } \Lambda_N = N$. Given a vector $\mathbf{y} \in \mathbb{R}^N$, the generalized interpolation problem can be formulated as an optimization problem to find

$$(28) \quad s_{\mathbf{y}; \Lambda_N} := \arg \min \left\{ \|s\|_{\mathcal{H}_\kappa(\mathbb{X})} : \lambda_k(s) = y_k \text{ for } 1 \leq k \leq N \right\}.$$

The solution to this problem turns out to be contained in a finite dimensional space

$$s_{\mathbf{y}; \Lambda_N} \in V_{\kappa; \Lambda_N} = \text{span} \left\{ \lambda_j^{(1)} \kappa(\cdot, \cdot) : \lambda_j \in \Lambda_N \right\}.$$

The coefficients can be computed by

$$s_{\mathbf{y}; \Lambda_N} = \sum_{j=1}^N \alpha_j \lambda_j^{(1)} \kappa(\cdot, \cdot), \text{ where } \begin{pmatrix} \lambda_1^{(1)} \lambda_1^{(2)} \kappa(\cdot, \cdot) & \dots & \lambda_1^{(1)} \lambda_N^{(2)} \kappa(\cdot, \cdot) \\ \vdots & \ddots & \vdots \\ \lambda_N^{(1)} \lambda_1^{(2)} \kappa(\cdot, \cdot) & \dots & \lambda_N^{(1)} \lambda_N^{(2)} \kappa(\cdot, \cdot) \end{pmatrix} \begin{pmatrix} \alpha_1 \\ \vdots \\ \alpha_N \end{pmatrix} = \begin{pmatrix} y_1 \\ \vdots \\ y_N \end{pmatrix}.$$

In order to apply this to a kernel-based collocation method for the numerical solution to (19), we consider a union of two collocation sets $Z_L = Z_{L_1} \cup Z_{L_2}$ with $Z_{L_1} := \{z_1, \dots, z_{L_1}\} \subset D$ and $Z_{L_2} := \{z_{L_1+1}, \dots, z_{L_2}\} \subset \partial D$. We define

$$\Lambda_L := \Lambda_{L_1} \cup \Lambda_{L_2} = \{\delta_{z_j} \circ L_c : z_j \in Z_{L_1}\} \cup \{\delta_{z_j} : z_j \in Z_{L_2}\} \subset W_2^{-\tau}(\mathbb{R}^d).$$

The data vector is given by

$$(29) \quad \mathbf{y} = (g(z_1), \dots, g(z_{L_1}), 0, \dots, 0)^\top \in \mathbb{R}^L = \mathbb{R}^{L_1+L_2}.$$

Moreover, we assume that there are open sets $V_k \subset \mathbb{R}^d$ such that

$$\partial D \subset \bigcup_{k=1}^{n_{\partial D}} V_k, \quad C^{\tau+2} \ni \phi_k : B(\mathbf{0}, 1) \xrightarrow{\simeq} V_k$$

and we denote by w_k a subordinate partition of unity which allows to write for $1 \leq p < \infty$

$$\|u\|_{W_p^\tau(\partial D)}^p := \sum_{k=1}^{n_{\partial D}} \|(u \cdot w_k) \circ \phi_k\|_{W_p^\tau(B)}^p$$

and the usual modifications for $p = \infty$. The discretization parameters are defined for a discrete set of point $P = \{p_1, \dots, p_{\#P}\} \subset A \subset \mathbb{R}^d$ as in (15), i.e.,

$$h_{P,A} := \sup_{\mathbf{a} \in A} \min_{\mathbf{p} \in P} \|\mathbf{a} - \mathbf{p}\|_2.$$

Hence, we can directly use this definition for $h_{Z_{L_1}, D}$. For the points on the boundary, we define

$$h_{Z_{L_2}, \partial D} := \max_{k=1}^{n_{\partial D}} h_{\phi_k^{-1}(Z_{L_2} \cap V_k), B(\mathbf{0}, 1)}.$$

Since the atlas is assumed to be fixed, the dependence on the particular atlas does not matter. Then, we can use [4, Corollary 3.13].

Proposition 4.1. *Let $\mathcal{L}_{D;\mathbf{c}}^{-1}(g, 0) \in W_2^\tau(D)$ be the solution to (19) and set $\hat{\mathcal{L}}_{D;\mathbf{c}_m;L}^{-1}(g, 0) := s_{\mathbf{y};\Lambda_L}$ with the notation from (28). Then, there is a constant $C > 0$ such that*

$$\begin{aligned} \left\| \mathcal{L}_{D;\mathbf{c}_m}^{-1}(g, 0) - \hat{\mathcal{L}}_{D;\mathbf{c}_m;L}^{-1}(g, 0) \right\|_{L^\infty(D)} &\leq C \left(h_{Z_{L_1}, D}^{\tau-2-\frac{d}{2}} + h_{Z_{L_2}, \partial D}^{\tau-\frac{1+(d-1)}{2}} \right) \|u\|_{W_2^\tau(D)} \\ (30) \qquad \qquad \qquad &\leq CC(A_{\min}, A_{\max}) \left(h_{Z_{L_1}, D}^{\tau-2-\frac{d}{2}} + h_{Z_{L_2}, \partial D}^{\tau-\frac{d}{2}} \right) \|g\|_{W_2^{\tau-2}(D)} \end{aligned}$$

for all $\mathbf{c} \in I^K$.

So far, (30) is an error estimate for the numerical solution of a differential equation.

Corollary 4.2. *We define a perturbed observation operator*

$$\hat{\mathcal{G}} : I^K \rightarrow \mathbb{R}^N, \quad \mathbf{c} \mapsto S_{X_N}(\mathcal{E}_D \circ \hat{\mathcal{L}}_{D;\mathbf{c}}^{-1}(g, 0)).$$

Then, we obtain an error bound for the difference of the perturbed observation operator and its unperturbed counterpart (22)

$$(31) \quad \left\| \hat{\mathcal{G}}(\mathbf{c}) - \mathcal{G}(\mathbf{c}) \right\|_{\ell_\infty(\mathbb{R}^N)} \leq \epsilon_L := CC(A_{\min}, A_{\max}) \left(h_{Z_{L_1}, D}^{\tau-2-\frac{d}{2}} + h_{Z_{L_2}, \partial D}^{\tau-\frac{d}{2}} \right) \|g\|_{W_2^{\tau-2}(D)}$$

for all $\mathbf{c} \in I^K$.

Proof. We directly observe for $\epsilon_{m,L}$ given in (23)

$$\|\epsilon_{m;L}\|_{\ell_\infty(X_N)} := \left\| S_{X_N} \circ \mathcal{E}_D \left(\mathcal{L}_{D;\mathbf{c}_m}^{-1}(g, 0) - \hat{\mathcal{L}}_{D;\mathbf{c}_m;L}^{-1}(g, 0) \right) \right\|_{\ell_\infty(X_N)} \leq \epsilon_L,$$

which is a uniform bound for $\mathbf{c}_m \in C_M \subset I^K$. \square

5. ERROR ANALYSIS FOR THE INVERSE PROBLEM

Here, we mainly follow [15]. The main goal of Bayesian inverse problems is to reconstruct the posterior distribution (25).

5.1. Reconstruction of observation operator. In this section, we consider the reconstruction of the observation operator (22) $\mathcal{G} : I^K \rightarrow \mathbb{R}^N$ given polluted evaluations at the points $C_M \subset I^K$. As a first step, we consider a fixed set of points and later study the incremental construction of the points sets. We recall the structure of the data \mathbf{y}_m for $1 \leq m \leq M$

$$\begin{aligned} \mathbf{y}_m &= \mathcal{G}(\mathbf{c}_m) = S_{X_N}(\mathcal{L}_{D;\mathbf{c}_m}^{-1}(g, 0)) \\ &= S_{X_N}(\hat{\mathcal{L}}_{D;\mathbf{c}_m;L}^{-1}(g, 0)) + S_{X_N}(\hat{\mathcal{L}}_{D;\mathbf{c}_m;L}^{-1}(g, 0) - \mathcal{L}_{D;\mathbf{c}_m}^{-1}(g, 0)) \\ &=: \bar{\mathbf{y}}_m + \epsilon_{m;L} = \bar{\mathbf{y}}_m, \end{aligned}$$

where $\bar{\mathbf{y}}_m = S_{X_N} \circ \mathcal{E}_D \circ \hat{\mathcal{L}}_{D;\mathbf{c}_m;L}^{-1}(g, 0)$ are the available observations which are polluted by the numerical model error given in (31).

Since $\mathcal{G} : I^K \rightarrow \mathbb{R}^N$, we write

$$\mathcal{G}(\mathbf{c}) = (\mathcal{G}_1(\mathbf{c}), \dots, \mathcal{G}_N(\mathbf{c}))^\top$$

and also

$$(32) \quad S_{C_M}(\mathcal{G}) = \begin{pmatrix} y_{1;1} & \cdots & y_{1;M} \\ \vdots & \ddots & \vdots \\ y_{1;N} & \cdots & y_{N;M} \end{pmatrix} \approx \bar{\mathbf{Y}} = \begin{pmatrix} \bar{y}_{1;1} & \cdots & \bar{y}_{1;M} \\ \vdots & \ddots & \vdots \\ \bar{y}_{1;N} & \cdots & \bar{y}_{N;M} \end{pmatrix} \in \mathbb{R}^{N \times M}$$

and treat each component separately. Since we have a (vector-valued) reconstruction problem for \mathcal{G} , we choose a regression approach instead of the interpolation approach of [15]. We assume that for some $[\sigma] > \frac{d}{2}$

$$\mathcal{G} \in W_2^\sigma(I^K)$$

and again $W_2^\sigma(\mathbb{R}^d)$ is an RKHS with reproducing kernel

$$K^{(\sigma)} : \mathbb{R}^K \times \mathbb{R}^K \rightarrow \mathbb{R}$$

which satisfies the conditions of (11) with $\tau = \sigma$ and $d = K$. We consider the functional

$$(33) \quad J_{\bar{\mathbf{y}}, \alpha, n} : W_2^\sigma(I^K) \rightarrow \mathbb{R}, \quad J_{\bar{\mathbf{y}}, \alpha, n}(s) := \sum_{m=1}^M (s(\mathbf{c}_m) - \bar{y}_{n;m})^2 + \alpha_n \|s\|_{W_2^\sigma(I^K)}^2.$$

The usual representer's theorem yields that an optimum exists and we define

$$(34) \quad s_n^* := \arg \min_{s \in W_2^\sigma(\mathbb{R}^d)} J_{\bar{\mathbf{y}}, \alpha, n}(s) \in \text{span} \left\{ K^{(\sigma)}(\cdot, \mathbf{c}_1), \dots, K^{(\sigma)}(\cdot, \mathbf{c}_M) \right\}.$$

Let $\hat{s}_n^* \in \mathbb{R}^M$ denote the coefficients of s_n^* , i.e.,

$$s_n^* = \sum_{m=1}^M \hat{s}_{n;m}^* K^{(\sigma)}(\cdot, \mathbf{c}_m).$$

We consider the symmetric positive definite matrix

$$\mathbf{K}_{C_M, C_M}^{(\sigma)} := \begin{pmatrix} K^{(\sigma)}(\mathbf{c}_1, \mathbf{c}_1) & \cdots & K^{(\sigma)}(\mathbf{c}_1, \mathbf{c}_M) \\ \vdots & \ddots & \vdots \\ K^{(\sigma)}(\mathbf{c}_M, \mathbf{c}_1) & \cdots & K^{(\sigma)}(\mathbf{c}_M, \mathbf{c}_M) \end{pmatrix} \in \mathbb{R}^{M \times M}.$$

The coefficients can be obtained by solving the linear system

$$(35) \quad \left(\mathbf{K}_{C_M, C_M}^{(\sigma)} + \alpha_n \mathbf{Id}_{M \times M} \right) \hat{\mathbf{s}}_n^* = \bar{\mathbf{y}}_{n;1:M} = \bar{\mathbf{Y}}^\top \mathbf{e}_n^{(N)},$$

where $\mathbf{e}_n^{(N)}$ is the n -th unit vector in \mathbb{R}^N . Note at this point, that the optimization problem (34) can be seen as a penalty methods for the optimization problem (28). For the error analysis, we follow [16].

Theorem 5.1. *For the error between the observation operator (22) and the reconstructed observation operator using the data (32) we obtain for $\mathcal{G}_{n;N,M} = s_n^*$ defined in (34) the error estimate*

$$\|\mathcal{G}_n - \mathcal{G}_{n;N,M}\|_{L_2(I^K)} \lesssim h_{C_M, I^K}^{\frac{K}{2}} \left(\sqrt{M} \epsilon_L + \sqrt{\alpha_n} \|\mathcal{G}_n\|_{W_2^\sigma(I^K)} \right) + h_{C_M, I^K}^\sigma \left(\frac{\sqrt{M}}{\sqrt{\alpha_n}} \epsilon_L + \|\mathcal{G}_n\|_{W_2^\sigma(I^K)} \right),$$

where $\mathcal{G}_{n;N,M}$ denotes the n -th component of $\mathcal{G}_{N,M}$. For quasi-uniform sets of points, the error estimate reduces to

(36)

$$\|\mathcal{G}_n - \mathcal{G}_{n;N,M}\|_{L_2(I^K)} \lesssim \left(1 + h_{C_M, I^K}^{\sigma - \frac{K}{2}} \frac{1}{\sqrt{\alpha_n}} \right) \epsilon_L + \left(h_{C_M, I^K}^{\frac{K}{2}} \sqrt{\alpha_n} + h_{C_M, I^K}^\sigma \right) \|\mathcal{G}_n\|_{W_2^\sigma(I^K)}.$$

Proof. From the definition of the optimization functional (33), we can conclude that

$$\begin{aligned} \sum_{m=1}^M (s_n^*(\mathbf{c}_m) - \mathcal{G}_n(\mathbf{c}_m))^2 &= \sum_{m=1}^M (s_n^*(\mathbf{c}_m) - \hat{y}_{n;m})^2 \\ &\leq 2 \sum_{m=1}^M (s_n^*(\mathbf{c}_m) - \bar{y}_{n;m})^2 + 2 \sum_{m=1}^M \left| \epsilon_{m;L} \cdot \mathbf{e}_n^{(N)} \right|^2 \leq 2J_{\bar{\mathbf{y}}, \alpha, n}(s_n^*) + 2M\epsilon_L^2 \\ (37) \quad &\leq 2J_{\bar{\mathbf{y}}, \alpha, n}(\mathcal{G}_n) + 2M\epsilon_L^2 \leq 2\alpha_n \|\mathcal{G}_n\|_{W_2^\sigma(I^K)}^2 + 2M\epsilon_L^2. \end{aligned}$$

Furthermore, we get

$$(38) \quad \alpha_n \|s_n^*\|_{W_2^\sigma(I^K)}^2 \leq J_{\bar{\mathbf{y}}, \alpha, n}(s_n^*) \leq 2\alpha_n \|\mathcal{G}_n\|_{W_2^\sigma(I^K)}^2 + 2M\epsilon_L^2.$$

Now, we can use [8, Theorem 3.5] to deduce that

$$\begin{aligned} \|s_n^* - \mathcal{G}_n\|_{L_2(I^K)} &\lesssim h_{C_M, I^K}^{\frac{K}{2}} \left(\sum_{m=1}^M (s_n^*(\mathbf{c}_m) - \mathcal{G}_n(\mathbf{c}_m))^2 \right)^{\frac{1}{2}} + h_{C_M, I^K}^\tau \|s_n^* - \mathcal{G}_n\|_{W_2^\tau(I^K)} \\ (39) \quad &\lesssim h_{C_M, I^K}^{\frac{K}{2}} \left(\sqrt{M} \epsilon_L + \sqrt{\alpha_n} \|\mathcal{G}_n\|_{W_2^\tau(I^K)} \right) + h_{C_M, I^K}^\tau \left(\frac{\sqrt{M}}{\sqrt{\alpha_n}} \epsilon_L + \|\mathcal{G}_n\|_{W_2^\tau(I^K)} \right), \end{aligned}$$

where we used (37) and (38) in the last step. If the set $C_M \subset I^K$ is quasi-uniformly distributed in the sense that there is a constant (uniform as $M \rightarrow \infty$)

$$q_{C_M} := \frac{1}{2} \min_{\substack{\mathbf{c}, \mathbf{d} \in C_M \\ \mathbf{c} \neq \mathbf{d}}} \|\mathbf{c} - \mathbf{d}\|_2 \leq h_{C_M, I^K} \leq Cq_{C_M},$$

we get $h_{C_M, I^K}^{\frac{K}{2}} \sqrt{M} \sim 1$. Hence, for quasi-uniform sets, the error bound (39) reduces to using $s_n^* = \mathcal{G}_{n;N,M}$

$$\|\mathcal{G}_n - \mathcal{G}_{n;N,M}\|_{L_2(I^K)} \lesssim \left(1 + h_{C_M, I^K}^{\sigma - \frac{K}{2}} \frac{1}{\sqrt{\alpha_n}} \right) \epsilon_L + \left(h_{C_M, I^K}^{\frac{K}{2}} \sqrt{\alpha_n} + h_{C_M, I^K}^\sigma \right) \|\mathcal{G}_n\|_{W_2^\sigma(I^K)}.$$

Hence, we obtain

$$\|\mathcal{G}_n - \mathcal{G}_{n;N,M}\|_{L_2(I^k; \mathbb{R}^N)}^2 = \sum_{n=1}^N \|s_n^* - \mathcal{G}_n\|_{L_2(I^k)}^2 \leq h_{C_M, I^k}^\sigma \|\mathcal{G}\|_{W_2^\sigma(I^k; \mathbb{R}^N)} + N\epsilon_L.$$

□

These estimates also motivate a choice for the regularization parameter.

Corollary 5.2. *For the choice we choose*

$$(40) \quad \alpha_n = \alpha^* = h_{C_M, I^k}^{2(\sigma - \frac{K}{2})} = h_{C_M, I^k}^{2\sigma - K},$$

we obtain

$$\|\mathcal{G}_n - \mathcal{G}_{n;N,M}\|_{L_2(I^k)} \lesssim h_{C_M, I^k}^\sigma \|\mathcal{G}_n\|_{W_2^\sigma(I^k)} + \epsilon_L,$$

where ϵ_L denotes the numerical error as given in (31).

Now we will consider the situation where we construct the set of points C_M in an incremental way.

5.2. Error estimates for incrementally constructed point sets. Having the linear system (35) in mind, we can not directly hope to get a low rank approximation of the matrix in that linear system. We will also use the choice for the regularization parameter (40). We consider as fine set $\Upsilon_{N_\infty} \subset \mathbb{X}$ in the general setting of Section 2 a set $C_M \subset I^k$.

Lemma 5.3. *Let $h_{C_M, I^k} \leq h_0$, where h_0 is be defined in (17) with $\mathbb{X} = I^k$. Then, for all $\epsilon \in (0, 1)$ there is a constant $c_\epsilon > 0$ such that for all sets Ξ_n created by the discrete power greedy method (12) with the kernel $K^{(\sigma)}$ which satisfies (11), we have that there is a constant $\theta = \sigma - \frac{d}{2}$ such that*

$$(41) \quad h_{\Xi_n, I^k} \leq c_\epsilon n^{-\frac{\theta}{\theta + \epsilon} \frac{1}{K}}.$$

Proof. The proof follows directly from Theorem 2.2. □

We obtain also a result about low rank matrix approximation as

Corollary 5.4. *There is a constant $c_\epsilon > 0$ such that for all sets Ξ_n created by the discrete power greedy method (12), the matrix \mathbf{A}_{Ξ_n} as in (9) satisfies the estimate*

$$\left\| \mathbf{K}_{C_M, C_M}^{(\sigma)} + \alpha_n \mathbf{Id}_{M \times M} - \mathbf{A}_{\Xi_n} \right\|_F \leq C_P M \left(n^{-\frac{\sigma}{K} + \frac{1}{2}} + n^{-\frac{\theta}{\theta + \epsilon} \frac{2\sigma - K}{K}} \right)$$

where α_n is chosen according to (40). If $\sigma \geq 3K$, we get

$$\left\| \mathbf{K}_{C_M, C_M}^{(\sigma)} + \alpha_n \mathbf{Id}_{M \times M} - \mathbf{A}_{\Xi_n} \right\|_F \leq C_P M n^{-\frac{\sigma}{K} + \frac{1}{2}}.$$

Proof. We can use Theorem 2.1 to obtain

$$\left\| \mathbf{K}_{C_M, C_M}^{(\sigma)} - \mathbf{A}_{\Xi_n} \right\|_F \leq C_P c_\epsilon M n^{-\frac{\sigma}{K} + \frac{1}{2}}.$$

Now, we observe

$$\|\alpha_n \mathbf{Id}_{M \times M}\|_F \leq \alpha_n M \leq M h_{C_M, I^k}^{2\sigma - K} \leq c_\epsilon M n^{-\frac{\theta}{\theta + \epsilon} \frac{2\sigma - K}{K}},$$

where we used (40) and (41). Hence the claim follows by the triangle inequality if we choose ϵ small enough such that

$$\frac{\theta}{\theta + \epsilon} 2 \frac{\sigma}{K} - 1 \geq \frac{\sigma}{K} + \frac{1}{2} \Leftrightarrow \left(\frac{2\theta}{\theta + \epsilon} - 1 \right) \frac{\sigma}{K} \geq \frac{3}{2}.$$

□

Finally, we get error estimates for the reconstruction of the observation operator using the incrementally constructed set of points.

Theorem 5.5. *The error between the observation operator (22) and the reconstructed observation operator using the data (32) we obtain for $\mathcal{G}_{k;N,\Xi_n} = s_k^*$ with the solution to the variational problem (34) the error estimate*

$$(42) \quad \|\mathcal{G}_k - \mathcal{G}_{k;N,\Xi_n}\|_{L_2(I^K)} \leq C\epsilon_L + c_\epsilon \left(\sqrt{n} n^{-\frac{\theta}{\theta+\epsilon} \frac{2\sigma-K}{2K}} + n^{-\frac{\theta}{\theta+\epsilon} \frac{\sigma}{K}} \right) \|\mathcal{G}_n\|_{W_2^\sigma(I^K)},$$

where we use the choice (40) for the regularization parameter and the point set Ξ_n is constructed by the discrete power greedy method (12)

Remark 5.6. *If the point set Ξ_n is quasi-uniform, the estimate (42) loses one order of convergence compared to an interpolation with a quasi-uniform set of points.*

Proof. The proof follows by combining Corollary 5.2 with Lemma 5.3. □

5.3. Error estimate for the posterior approximation. The error for the posterior distribution is measured in the *Hellinger distance*. Let ν_i , $i = 1, 2$ be two measures which are absolutely continuous with respect to the prior measure ν_0 . Then the Hellinger distance between ν_1 and ν_2 is given as

$$\text{dist}_{\text{hell}}(\nu_1, \nu_2) := \left(\frac{1}{2} \int_{L_2(I^K)} \left(\sqrt{\frac{d\nu_1}{d\nu_0}} - \sqrt{\frac{d\nu_2}{d\nu_0}} \right) d\nu_0 \right)^{\frac{1}{2}}.$$

We define an approximate posterior by (27)

$$\frac{d\nu^{\mathbf{y}}}{d\nu_0}(\mathbf{c}) \approx \frac{1}{\mathbb{E}_{\mu_0} [\exp(-\|\mathbf{y} - \mathcal{G}(\mathbf{c})\|_2)]} \exp(-\|\mathbf{y} - \mathcal{G}_{N,M}(\mathbf{c})\|_2).$$

We can use [15, Theorem 4.2] to obtain

$$(43) \quad \text{dist}_{\text{hell}}(\nu, \nu^{\mathbf{y}}) \leq C \|\mathcal{G} - \mathcal{G}_{:,M,n}\|_{L_2(I^K)},$$

where the posterior is defined in (25). Hence, we obtain directly the following theorem

Theorem 5.7.

$$(44) \quad \text{dist}_{\text{hell}}(\nu, \nu^{\mathbf{y}}) \leq C\epsilon_L + c_\epsilon \left(\sqrt{n} n^{-\frac{\theta}{\theta+\epsilon} \frac{2\sigma-K}{2K}} + n^{-\frac{\theta}{\theta+\epsilon} \frac{\sigma}{K}} \right) \|\mathcal{G}_n\|_{W_2^\sigma(I^K)}$$

6. CONCLUSION

We re-interpreted the covariance of the predictive distribution as reproducing kernel. This interpretation allowed to reformulate certain (simple instances) of techniques from low-rank matrix approximation as simple greedy methods in reproducing kernel Hilbert spaces. In particular, this allowed to use results from the kernel-based literature to derive error estimates for the *adaptive cross approximation* or the *pivoted Cholesky* method as matrix approximation for finitely smooth kernels. We applied the resulting iterative construction of sampling point sets to the inverse problem to reconstruct a parametrized diffusion coefficient by finitely many point evaluations of the solution to the diffusion equation with that coefficient. We explicitly took the error in the solution to the differential equation also into account.

ACKNOWLEDGEMENT

The author acknowledges partial support of the Deutsche Forschungsgemeinschaft (DFG) through the Sonderforschungsbereich 1060: The Mathematics of Emergent Effects.

LITERATUR

- [1] M. BEBENDORF, *Adaptive Cross Approximation of Multivariate Functions*, Constructive Approximation, 34 (2011), pp. 149–179.
- [2] J. COCKAYNE, C. OATES, T. SULLIVAN, AND M. GIROLAMI, *Probabilistic Meshless Methods for Partial Differential Equations and Bayesian Inverse Problems*, May 2016. available at <https://arxiv.org/pdf/1605.07811v1.pdf>.
- [3] L. EVANS, *Partial Differential Equations*, Graduate studies in mathematics, American Mathematical Society, 2010.
- [4] P. GIESL AND H. WENDLAND, *Meshless collocation: Error estimates with application to dynamical systems*, SIAM Journal on Numerical Analysis, 45 (2007), pp. 1723–1741.
- [5] M. GRIEBEL AND C. RIEGER, *Reproducing Kernel Hilbert Spaces for Parametric Partial Differential Equations*, SIAM/ASA Journal on Uncertainty Quantification, 5 (2017), pp. 111–137.
- [6] H. HARBRECHT, M. PETERS, AND R. SCHNEIDER, *On the low-rank approximation by the pivoted Cholesky decomposition*, Appl. Numer. Math., 4 (2012), pp. 428–440.
- [7] R. KEMPF, H. WENDLAND, AND C. RIEGER, *Kernel-based Reconstructions for Parametric PDEs*. Available as INS Preprint No. 1804., 2018.
- [8] W. MADYCH, *An estimate for multivariate interpolation ii*, Journal of Approximation Theory, 142 (2006), pp. 116 – 128.
- [9] S. D. MARCHI, R. SCHABACK, AND H. WENDLAND, *Near-optimal data-independent point locations for radial basis function interpolation*, Advances in Computational Mathematics, 23 (2005), pp. 317–330.
- [10] M. MOUATTAMID AND R. SCHABACK, *Recursive Kernels*, Analysis in Theory and Applications, 25 (2009), p. 301.
- [11] S. MÜLLER AND R. SCHABACK, *A Newton basis for kernel spaces*, J. Approx. Theory, 161 (2009), pp. 645–655.
- [12] M. PAZOUKI AND R. SCHABACK, *Bases for kernel-based spaces*, Journal of Computational and Applied Mathematics, 236 (2011), pp. 575 – 588. International Workshop on Multivariate Approximation and Interpolation with Applications (MAIA 2010).
- [13] C. E. RASMUSSEN AND C. K. I. WILLIAMS, *Gaussian Processes for Machine Learning*, The MIT Press, 2006.
- [14] G. SANTIN AND B. HAASDONK, *Convergence rate of the data-independent P-greedy algorithm in kernel-based approximation*, Dolomites Research Notes on Approximation, 10 (2017), pp. 68–78.
- [15] A. M. STUART AND A. L. TECKENTRUP, *Posterior Consistency for Gaussian Process Approximations of Bayesian Posterior Distributions*, Mathematics of Computation, 87 (2018), pp. 721–753. available at: <https://arxiv.org/abs/1603.02004>.

- [16] H. WENDLAND AND C. RIEGER, *Approximate Interpolation with Applications to Selecting Smoothing Parameters*, *Numerische Mathematik*, 101 (2005), pp. 729–748.