# Deterministic error bounds for kernel-based learning techniques under bounded noise

**Emilio T. Maddalena**[1]                                        EMILIO.MADDALENA@EPFL.CH
**Paul Scharnhorst**[1,2]                                         PAUL.SCHARNHORST@EPFL.CH
**Colin N. Jones**[1]                                             COLIN.JONES@EPFL.CH

[1]*École Polytechnique Fédérale de Lausanne, Lausanne, Switzerland*
[2]*Swiss Center for Electronics and Microtechnology, Neuchâtel, Switzerland*

## Abstract

We consider the problem of reconstructing a function from a finite set of noise-corrupted samples. Two kernel algorithms are analyzed, namely kernel ridge regression and $\varepsilon$-support vector regression. By assuming the ground-truth function belongs to the reproducing kernel Hilbert space of the chosen kernel, and the measurement noise affecting the dataset is bounded, we adopt an approximation theory viewpoint to establish *deterministic*, finite-sample error bounds for the two models. Finally, we discuss their connection with Gaussian processes and two numerical examples are provided. In establishing our inequalities, we hope to help bring the fields of non-parametric kernel learning and system identification for robust control closer to each other.

**Keywords:** Deterministic error bounds; Generalization error; Kernel ridge regression; Support vector machines.

## 1. Introduction

As opposed to classical system identification techniques, where the structure of a finite-dimensional model is chosen *a priori*, kernel-based methodologies deal with possibly infinite-dimensional hypothesis spaces. In the latter, the number of parameters to be determined is not fixed, but depends on the number of available data-points. This non-parametric approach to building models is popular in many disciplines, usually in the form of Gaussian processes (GPs)—whose means are weighted sums of kernels—or plain radial basis functions (RBFs) (Jakhetiya et al., 2020; Moriconi et al., 2020; Singh et al., 2019). Among systems and control researchers, kernel methods have also been studied with the aim of adapting and improving existing tools (see Pillonetto et al. (2014); Chiuso and Pillonetto (2019); Ljung et al. (2020) for recent reviews). For instance, an in-depth analysis of linear system identification through stable kernels, which encode the asymptotic decay of the plant impulse response, was carried out in Pillonetto and De Nicolao (2010); Carli et al. (2016); Chen (2018); see also the work Lataire and Chen (2016) for a frequency domain perspective of the same problem, and Blanken and Oomen (2020) for the case of non-causal dynamical systems. Kernel-based identification of Hammerstein and Wiener systems, i.e. linear systems in cascade with a static nonlinearity, was studied from a similar point of view in Risuleo et al. (2017, 2019).

When nonlinear dynamics are considered in their full generality as in this work, a commonly adopted approach is to directly assume the availability of pairs of current and next states —that is, state-space models (Koller et al., 2018; Bradford et al., 2020; Umlauft and Hirche, 2020). The goal then is to reconstruct the complete vector field rather than the time-response of the unknown system

(see Pillonetto et al. (2011) for an exception to this statement), a function reconstruction problem remarkably similar to the ones found in machine learning. Certainly, embedding prior knowledge regarding for instance stability can be quite challenging in this context (Umlauft and Hirche, 2020).

Exploiting machine learning techniques to model and subsequently control physical systems requires caution, especially in safety-critical applications. The fact that GPs provide users with not only nominal predictions, but also confidence intervals, is arguably one of the main reasons for their popularity (Bradford et al., 2019; Binder et al., 2019; Hewing et al., 2020). This probabilistic uncertainty measure can then be used to assess the risk associated with actions, thus enabling the use of stochastic analysis techniques (Polymenakos et al., 2019; Jackson et al., 2020). If a non-probabilistic viewpoint is taken instead, it is possible to derive hard prediction-error bounds for the learned models such that the unknown ground-truth cannot lie outside the established 'prediction envelope'. These types of guarantees are widely known in the field of approximation theory (Schaback, 2000; Wendland, 2004; Fasshauer, 2011; Wang et al., 2019), where authors often consider the problem of interpolating noise-free observations. Whereas having access to perfect measurements might be common in domains such as computer graphics, it is definitely not the case in power networks, robotics, building automation systems, etc. Our goal herein is to extend such theory to the scenario where the designer is confronted with unknown but bounded noise, hence providing new tools for deterministic safety certification and robust control. Besides the latter application domain, the developed theory could also be of use in branches of natural sciences where the bounded noise assumption is regarded as more adequate than the Gaussian one (d'Onofrio, 2013).

**Contributions:** In this note, we study the problem of learning an unknown real-valued function $f$ from a set of evaluations corrupted by noise. To achieve this goal, we employ two distinct non-parametric kernel techniques, namely the popular kernel ridge regression (KRR), and $\varepsilon$-support vector regression (SVR) (Schölkopf et al., 2018). Two main assumptions are made[1]: firstly, the measurement noise is bounded by a known finite quantity; secondly, given a kernel $k$, the unknown $f$ belongs to the reproducing kernel Hilbert space (RKHS) $\mathcal{H}$ of $k$, and an upper bound for its norm $\|f\|_{\mathcal{H}}$ is known. The same RKHS assumptions were also exploited in the works Srinivas et al. (2012); Koller et al. (2018); Umlauft and Hirche (2019); Hashimoto et al. (2020). We then establish deterministic prediction-error bounds—otherwise known as risk or generalization error bounds—which extend classical results valid for noise-free interpolants only. Our final expressions are moreover given in closed-form, requiring only the solution of a simple box-constrained quadratic program in the KRR case. Finally, two numerical experiments are presented and a comparison with an existing alternative is provided.

**Notation:** $\mathbb{N} := \{1, 2, \dots\}$, $\mathbb{R}^n$ is the $n$-dimensional Euclidean space, and $\mathbb{R}^n_{\geq 0}$ ($\mathbb{R}^n_{>0}$) its positive (strictly positive) orthant. Given a matrix $A$, $A^\top$ denotes its transpose, and $A \succ 0$ indicates that $A$ is positive-definite. $I$ and $\mathbf{0}$ are respectively the identity and the zero matrices of appropriate sizes. Unless otherwise specified, $f$ denotes a map and $f(x)$, its particular evaluation at a point $x$. Given two vectors $a$ and $b$, the inequality $a \leq b$ and the absolute value $|a|$ are to be read element-wise. $\|a\|_1$ and $\|a\|_2$ will be used respectively for the 1-norm and 2-norm of a vector, whereas $\|f\|_{\mathcal{H}}$, for the RKHS functional norm.

---

1. Note that these are fundamentally different compared to the ones made in the Gaussian processes setting. For more information, please refer to the discussion in Subsection 4.3.

## 2. Problem setting

Consider the problem of learning an unknown map $f : \mathcal{X} \to \mathbb{R}$, referred to as the ground-truth or target function. Herein $\mathcal{X} \subset \mathbb{R}^m$ is assumed to be compact. In order to reconstruct $f$, we collect a finite dataset

$$D = \{(x_n, y_n) \,|\, n = 1, \ldots, N\} \tag{1}$$

composed of sites $x_n$ and noisy evaluations of the ground-truth

$$y_n = f(x_n) + \delta_n, \ n = 1, \ldots, N \tag{2}$$

**Assumption 1** *The data sites $X = \{x_1, \ldots, x_N\}$ in $D$ are pairwise distinct.*

**Assumption 2** *The measurement noise $\delta = \begin{bmatrix} \delta_1 & \ldots & \delta_N \end{bmatrix}^\top$ is bounded $|\delta| \leq \bar{\delta}$, with $\bar{\delta} \in \mathbb{R}^N_{\geq 0}$.*

Working with bounded uncertainties is at the core of robust analysis and control. Similar assumptions were also made in recent learning-based techniques (Rosolia et al., 2017; Novara et al., 2019; Manzano et al., 2020). Next, we review basic definitions and results from RKHS theory, which are used as a starting point for Section 3. The reader is referred to Manton and Amblard (2015) for a more in-depth discussion on the topic.

### 2.1. Kernels and their RKHSs

We call a symmetric real-valued map $k : \mathcal{X} \times \mathcal{X} \to \mathbb{R}$ a kernel, and assume $k$ is a positive-definite (PD) function according to the definition that follows.

**Definition 1** *A continuous function $k : \mathcal{X} \times \mathcal{X} \to \mathbb{R}$ is called positive-definite if for any set of pairwise-distinct sites $X = \{x_1, \ldots, x_N\}$, with an arbitrary $N \in \mathbb{N}$, it holds that*

$$\sum_{i=1}^{N} \sum_{j=1}^{N} \alpha_i \alpha_j k(x_i, x_j) > 0$$

*for any set of weighting constants $\alpha_1, \ldots, \alpha_N \in \mathbb{R}\backslash\{0\}$.*

Even though limiting our scope to positive-definite functions excludes certain kernels, this class encompasses powerful alternatives such as the squared-exponential, the inverse multiquadrics and the truncated power function (Wendland, 2004). The first of the three is known to have the universal approximation property whenever $\mathcal{X}$ is compact (Micchelli et al., 2006).

**Remark 1** *For clarity purposes, we denote by $K \in \mathbb{R}^{N \times N}$ the constant matrix that has kernel evaluations as elements, i.e., $k(x_i, x_j)$ at its ith row and jth column for $x_i, x_j \in X$. Moreover, $K_{Xx} : \mathcal{X} \to \mathbb{R}^N$ denotes the column vector function $x \mapsto \begin{bmatrix} k(x_1, x) & \ldots & k(x_N, x) \end{bmatrix}^\top$, and $K_{xX}$ simply represents its transpose.*

Given a kernel $k$, we denote the associated uniquely determined reproducing kernel Hilbert space (RKHS) by $\mathcal{H}$. The one-to-one correspondence between $k$ and $\mathcal{H}$ is guaranteed by the Moore–Aronszajn theorem (Aronszajn, 1950). Each element $g \in \mathcal{H}$ is a map from $\mathcal{X}$ to $\mathbb{R}$ assuming the form of a weighted sum of kernels $g = \sum_{i \in \Omega_g} \alpha_i k(x_i, \cdot)$, where the index set $\Omega_g$ of $g$ can also be countably infinite, in which case the limit interpretation of the series applies. $\mathcal{H}$ is equipped with the inner product $\langle g, f \rangle_{\mathcal{H}} = \sum_{i \in \Omega_g} \sum_{j \in \Omega_f} \alpha_i \beta_j k(x_i, x_j)$ and the induced norm is $\|g\|_{\mathcal{H}} := \sqrt{\langle g, g \rangle_{\mathcal{H}}}$. Fixing any $x$ in $\mathcal{X}$, the corresponding evaluation functional $l_x : \mathcal{H} \to \mathbb{R}$ is bounded and takes any $g \in \mathcal{H}$ to its image, i.e., $l_x(g) = g(x) = \langle g, k(x, \cdot) \rangle_{\mathcal{H}}$, being linked to the reproducing property. Suppose that $g$ has a finite expansion in terms of $N_g$ kernel functions. Due to the reproducing and basic inner product properties, it holds that

$$\|g\|_{\mathcal{H}}^2 = \left\langle \sum_{i=1}^{N_g} \alpha_i k(x_i, \cdot), \sum_{i=1}^{N_g} \alpha_i k(x_i, \cdot) \right\rangle_{\mathcal{H}} \tag{3}$$

$$= \sum_{i=1}^{N_g} \sum_{j=1}^{N_g} \alpha_i \alpha_j k(x_i, x_j) \tag{4}$$

$$= \alpha^\top K \alpha \tag{5}$$

where $\alpha := \begin{bmatrix} \alpha_1 & \ldots & \alpha_{N_g} \end{bmatrix}^\top$ gathers all weights.

**Assumption 3** *Given a kernel $k$, we assume that the ground-truth $f$ belongs to its RKHS, $\mathcal{H}$. Additionally, an upper bound for its norm $\|f\|_{\mathcal{H}} \leq \Gamma$ is available.*

Establishing any form of guarantee is clearly impossible if no assumptions are made on the unknown map $f$. The availability of an upper bound $\Gamma$ is also assumed in the works Srinivas et al. (2012); Koller et al. (2018); Hashimoto et al. (2020). It is beyond the scope of this paper to present various ways to construct bounds for $\|f\|_{\mathcal{H}}$. We nevertheless illustrate how this quantity can be estimated from perfect evaluations of $f$ in Section 5. The following function is defined as a last introductory step, which will later play an important role in our error estimates.

**Definition 2** *The power function is the real-valued map $P_X(x) = \sqrt{k(x, x) - K_{xX} K^{-1} K_{Xx}}$.*

Throughout most of the document, $D$ is assumed to be fixed. For this reason, we drop the dependence that $P_X(x)$ has on the data sites to ease the notation, writing simply $P(x)$. Two main properties of this function will be exploited herein:

$$P(x) \geq 0, \text{ for any } x \in \mathcal{X}$$

$$P(x_n) = 0, \text{ for any } x_n \in X$$

which follow from rewriting the power function in a Lagrange form as shown in (Wendland, 2004, Sec. 11.1).

## 3. Crafting models

We restrict our attention to models $s : \mathcal{X} \to \mathbb{R}$ built as a weighted sum of kernels that are centered at the data locations

$$s(x) = \sum_{n=1}^{N} \alpha_n k(x_n, x) = \alpha^\top K_{Xx}. \tag{6}$$

Solutions to a number of optimal fitting problems have this form as discussed next. Since the number of functions $k$ and their centers have already been defined, constructing a model is equivalent to deciding the $\alpha$ coefficients.

### 3.1. The noise-free case

In the absence of noise ($\bar{\delta} = \mathbf{0}$), the labels in $D$ perfectly represent $f$. We can then solve the approximation problem by finding an $s \in \mathcal{H}$ such that the evaluations $s(x_i)$ match $f(x_i) =: f_{x_i}$ for all points in $D$. This can be posed as the variational problem

$$\bar{s} = \underset{s \in \mathcal{H}}{\arg\min} \quad \|s\|_{\mathcal{H}}^2 \tag{7a}$$

$$\text{subj. to} \quad s(x_n) = f_{x_n} \tag{7b}$$
$$\forall n = 1, \dots, N$$

in which the objective favors low-complexity solutions, measured by the function space norm $\|\cdot\|_{\mathcal{H}}$.

Thanks to the optimal recovery property (see (Wendland, 2004, Thm 13.2) or (Kanagawa et al., 2018, Thm. 3.5)), it is known that out of all elements $s \in \mathcal{H}$ capable of interpolating the dataset, a minimizer for the above problem exists and assumes the form (6). The solution $\bar{s}$ can be therefore found by simply solving the linear system of equations $K\alpha = f_X$ for $\alpha$, where $f_X = \begin{bmatrix} f(x_1) & \dots & f(x_N) \end{bmatrix}^\top$. Given the PD property of the kernel $k$ and Assumption 1, $K$ is positive-definite and hence invertible. Therefore, $\alpha = K^{-1} f_X$ and the unique optimizer of (7) is

$$\bar{s}(x) = f_X^\top K^{-1} K_{Xx} \tag{8}$$

Because of (5), we see that its norm can be expressed in terms of the data values as $\|\bar{s}\|_{\mathcal{H}}^2 = f_X^\top K^{-1} f_X$.

**Remark 2** *It holds that $\|\bar{s}\|_{\mathcal{H}} \leq \|f\|_{\mathcal{H}}$ independently of the number of samples in $D$. This stems from $f$ being the solution to (7) when the equality constraints are imposed for all $x \in \mathcal{X}$.*

A first inequality can be obtained for the model $\bar{s}$ with the aid of the previous remark. This is a known but not very disseminated result, which tightens the more widely spread bound (Wendland, 2004, Thm. 11.4). The proof we give here is important to help build an intuition on how the RKHS norm measures the complexity of a function.

**Proposition 1** *Assume that the dataset $D$ is not affected by noise, i.e., $\bar{\delta} = \mathbf{0}$ and $y = f_X$. Under Assumptions 1 to 3, the interpolating model $\bar{s}$ admits the error bound*

$$|\bar{s}(x) - f(x)| \leq P(x) \sqrt{\Gamma^2 - \|\bar{s}\|_{\mathcal{H}}^2} \tag{9}$$

*for any $x \in \mathcal{X}$, where $f$ is the unknown ground-truth and $\|\bar{s}\|_{\mathcal{H}}^2 = f_X^\top K^{-1} f_X$.*

**Proof** Let $x \in \mathcal{X}$ be a fixed query point, which is not in $D$. Denote by $\bar{s}_+$ the function of the form (6) interpolating all known points $f_X$ in $D$ and the unknown value $f_x := f(x)$. We then have

$$
\begin{aligned}
\|\bar{s}_+\|_{\mathcal{H}}^2 &= \begin{bmatrix} f_X \\ f_x \end{bmatrix}^\top \begin{bmatrix} K & K_{Xx} \\ K_{xX} & K_{xx} \end{bmatrix}^{-1} \begin{bmatrix} f_X \\ f_x \end{bmatrix} \\
&= \begin{bmatrix} f_X \\ f_x \end{bmatrix}^\top \begin{bmatrix} K^{-1} & \mathbf{0} \\ \mathbf{0} & 0 \end{bmatrix} \begin{bmatrix} f_X \\ f_x \end{bmatrix} + P^{-2}(x) \begin{bmatrix} f_X \\ f_x \end{bmatrix}^\top \begin{bmatrix} K^{-1} K_{Xx} \\ -1 \end{bmatrix} \begin{bmatrix} K^{-1} K_{Xx} \\ -1 \end{bmatrix}^\top \begin{bmatrix} f_X \\ f_x \end{bmatrix} \\
&= \|\bar{s}\|_{\mathcal{H}}^2 + P^{-2}(x)(\bar{s}(x) - f_x)^2 \\
&\leq \Gamma^2
\end{aligned}
$$

where the second equality follows from the matrix inversion lemma, and the inequality follows from Remark 2. Finally, the last two lines imply $|\bar{s}(x) - f(x)| \leq P(x)\sqrt{\Gamma^2 - \|\bar{s}\|_{\mathcal{H}}^2}$.

If, on the other hand, the query point $x$ belongs to the dataset $D$, the bound evaluates to zero and thus it holds tightly. ∎

One can also arrive at (9) by starting from the inequality $|\bar{s}(x) - f(x)| \leq P(x)\|f - \bar{s}\|_{\mathcal{H}}$ (Fasshauer, 2011, Eq. 9) and noting that $\|f - \bar{s}\|_{\mathcal{H}}^2 = \|f\|_{\mathcal{H}}^2 - 2\langle \bar{s}, f \rangle_{\mathcal{H}} + \|\bar{s}\|_{\mathcal{H}}^2 = \|f\|_{\mathcal{H}}^2 - 2\|\bar{s}\|_{\mathcal{H}}^2 + \|\bar{s}\|_{\mathcal{H}}^2 \leq \Gamma^2 - \|\bar{s}\|_{\mathcal{H}}^2$. Nevertheless, the proof we provided gives more insight on how the norm of a model $\bar{s}$ might grow after the addition of a new data-point. More specifically, this only happens if the new value $f_x$ differs from what the model was previously predicting $\bar{s}(x)$.

Through Proposition 1, evaluations of $f$ for every $x \in \mathcal{X}$ can be bounded according to $f^{\min}(x) \leq f(x) \leq f^{\max}(x)$ with $f^{\min}(x) = \bar{s}(x) - P(x)\sqrt{\Gamma^2 - \|\bar{s}\|_{\mathcal{H}}^2}$ and $f^{\max}(x) = \bar{s}(x) + P(x)\sqrt{\Gamma^2 - \|\bar{s}\|_{\mathcal{H}}^2}$. Proposition 2 below establishes that the interval containing the ground-truth is non-growing after the addition of a new data-point.

**Proposition 2** *Let $x$ be any fixed query point in the domain $\mathcal{X}$. Let $Z$ be an augmented set of distinct data-sites, i.e., $Z = X \cup \{z\}$, $z \in \mathcal{X}, z \notin X$. Then we have*

$$
[f_Z^{\min}(x), f_Z^{\max}(x)] \subseteq [f_X^{\min}(x), f_X^{\max}(x)] \tag{11}
$$

**Proof** See Appendix A. ∎

Intuitively, this favorable property states that augmenting the dataset $D$ with any new pair $(x, y) \in \mathcal{X} \times \mathbb{R}$ (while still satisfying Assumption 1) either preserves or sharpens the inequality (9). This holds everywhere in the domain.

### 3.2. Kernel ridge regression analysis

To tackle the approximation problem in the presence of measurement noise, a compromise between fitting the data and rejecting uninformative fluctuations has to be found. One of the most standard tools used to achieve this balance is kernel ridge regression (KRR), in which the unconstrained problem

$$
s^* = \arg\min_{s \in \mathcal{H}} \frac{1}{N} \sum_{n=1}^{N} (y_n - s(x_n))^2 + \lambda \|s\|_{\mathcal{H}}^2 \tag{12}
$$

is solved, yielding the KRR model $s^*$. In (12), the regularization weight $\lambda \in \mathbb{R}_{\geq 0}$ dictates the aforementioned balance: $\lambda = 0$ produces a plain interpolant due to our assumptions on $k$ and $D$, whereas increasing values of $\lambda$ lead to a surrogate model that fits the dataset while avoiding abrupt variations. See Scandella et al. (2020) for a recent study on the numerical properties of the problem.

Thanks to the representer theorem, this infinite-dimensional functional problem over $\mathcal{H}$ can be recast as an equivalent finite-dimensional one (see the pivotal work Schölkopf et al. (2001), and Diwale and Jones (2018) for a recent generalization). A closed-form solution for (12) can then be obtained through this reformulation and is given by (see e.g. (Kanagawa et al., 2018, Thm. 3.4))

$$s^*(x) = \alpha^{*\top} K_{Xx} \tag{13}$$

with optimal weights $\alpha^* = (K + N\lambda I)^{-1} y$. Let $c$ denote the vector of values attained by $s^*$ at the data locations $X$, i.e., $c_n := s^*(x_n), n = 1, \ldots, N$. The regressor then satisfies $K\alpha^* = c \Rightarrow \alpha^* = K^{-1} c$. From the latter and (5), the norm can be also expressed in the convenient quadratic form $\|s^*\|_{\mathcal{H}}^2 = c^\top K^{-1} c$, where $c = K(K + N\lambda I)^{-1} y$.

Before establishing the KRR prediction-error bound, we first need to analyze how noise can perturb the norm of an interpolant. To this end, consider the next result.

**Lemma 1** *Let Assumptions 1 and 2 hold. Let moreover $\bar{s}(x) = f_X^\top K^{-1} K_{Xx}$ be the model interpolating the noise-free values $f_X$, and $\tilde{s}(x) = y^\top K^{-1} K_{Xx}$ the model interpolating the noisy values $y$. Then*

$$\nabla \leq \|\tilde{s}\|_{\mathcal{H}}^2 - \|\bar{s}\|_{\mathcal{H}}^2 \leq \Delta \tag{14}$$

*where $\Delta$ denotes the maximum and $\nabla$ the minimum of $(-\delta^\top K^{-1} \delta + 2 y^\top K^{-1} \delta)$ subject to $|\delta| \leq \bar{\delta}$.*

**Proof** It follows from expanding $\|\tilde{s}\|_{\mathcal{H}}^2$ as $\|\bar{s}\|_{\mathcal{H}}^2$ plus a perturbation term, and recalling the definitions of $\Delta$ and $\nabla$. ■

Whereas calculating $\Delta$ amounts to solving a convex optimization problem since it is the maximum of a strictly concave function, evaluating $\nabla$ is not as straightforward. Still, the quantity $\nabla$ is not employed in our expressions. Bounding the error associated with the KRR predictions is then possible through the following inequality.

**Theorem 1** *Let $N$ be the number of data-points, $\bar{\delta} \in \mathbb{R}_{>0}^N$ the noise bound, and $\lambda$ the regularization constant. Under Assumptions 1 to 3, the KRR model $s^*$ admits the error bound*

$$|s^*(x) - f(x)| \leq P(x) \sqrt{\Gamma^2 + \Delta - \|\tilde{s}\|_{\mathcal{H}}^2} + \bar{\delta}^\top |K^{-1} K_{Xx}| + \left| y^\top \left( K + \frac{1}{N\lambda} KK \right)^{-1} K_{Xx} \right| \tag{15}$$

*for any $x \in \mathcal{X}$, where $f$ is the unknown ground-truth, $\Delta = \max_{|\delta| \leq \bar{\delta}} (-\delta^\top K^{-1} \delta + 2 y^\top K^{-1} \delta)$, and $\|\tilde{s}\|_{\mathcal{H}}^2 = y^\top K^{-1} y$.*

**Proof** Recall that $\bar{s}(x) = f_X^\top K^{-1} K_{Xx}$ and $y = f_X + \delta$. Predictions given by $s^*(x)$ can be decomposed as

$$s^*(x) = y^\top (K + N\lambda I)^{-1} K_{Xx} \tag{16a}$$

$$= f_X^\top (K + N\lambda I)^{-1} K_{Xx} + \delta^\top (K + N\lambda I)^{-1} K_{Xx} \tag{16b}$$

$$= f_X^\top K^{-1} K_{Xx} - f_X^\top (K + \frac{1}{N\lambda} KK)^{-1} K_{Xx} + \delta^\top (K + N\lambda I)^{-1} K_{Xx} \tag{16c}$$

$$= f_X^\top K^{-1} K_{Xx} - y^\top (K + \frac{1}{N\lambda} KK)^{-1} K_{Xx} + \delta^\top \left[ (K + \frac{1}{N\lambda} KK)^{-1} + (K + N\lambda I)^{-1} \right] K_{Xx} \tag{16d}$$

$$= \bar{s}(x) - y^\top (K + \frac{1}{N\lambda} KK)^{-1} K_{Xx} + \delta^\top K^{-1} K_{Xx} \tag{16e}$$

where (16c) and (16e) both follow from Woodbury's matrix identity. For compactness, let $Q = (K + \frac{1}{N\lambda} KK)$. The error norm can therefore be upper bounded by

$$|s^*(x) - f(x)| = |\bar{s}(x) - f(x) + \delta^\top K^{-1} K_{Xx} - y^\top Q^{-1} K_{Xx}| \tag{17a}$$

$$\leq |\bar{s}(x) - f(x)| + |\delta^\top K^{-1} K_{Xx} - y^\top Q^{-1} K_{Xx}| \tag{17b}$$

$$\leq |\bar{s}(x) - f(x)| + \bar{\delta}^\top |K^{-1} K_{Xx}| + |y^\top Q^{-1} K_{Xx}| \tag{17c}$$

$$\leq P(x)\sqrt{\Gamma^2 - \|\bar{s}\|_{\mathcal{H}}^2} + \bar{\delta}^\top |K^{-1} K_{Xx}| + |y^\top Q^{-1} K_{Xx}| \tag{17d}$$

in which the triangle inequality was used to obtain (17b). Note that (17c) *tightly* bounds (17b), i.e., $\exists \delta : |\delta| \leq \bar{\delta}$ such that both expressions are equal. Finally, (17d) is due to Proposition 1. Thanks to Lemma 1, we know that $\|\tilde{s}\|_{\mathcal{H}}^2 - \Delta \leq \|\bar{s}\|_{\mathcal{H}}$, concluding the proof. ∎

Let us have a closer look at the bound (15). Firstly, it is consistent with the interpolating noise-free case, i.e., if $\bar{\delta} = \mathbf{0}$ and $\lambda \to 0$, (15) converges to (9). Secondly, the constant $\Delta$ was introduced since we do not have access to $\|\bar{s}\|_{\mathcal{H}}^2$. Calculating it requires solving a box-constrained quadratic program (QP) for it is the maximization of a concave objective. The maximum is moreover *independent* of the query point $x$. As a compelling alternative, note that $P(x)\Gamma$ can be used as a replacement for $P(x)\sqrt{\Gamma^2 + \Delta - \|\tilde{s}\|_{\mathcal{H}}^2}$ in (15), being a simple consequence of (9). By doing that, we observed experimentally that little conservativeness is introduced, while avoiding the need to compute $\Delta$.

Approaching the problem from another perspective, we analyze now a technique that enjoys a convenient low-norm solution, which in turn leads to simplified error bounds.

### 3.3. $\varepsilon$-Support vector regression analysis

We now aim at finding a minimum norm solution to the approximation problem under noise, while also fully exploiting the boundedness of $\delta$. This can be readily formulated as

$$s^\star = \arg\min_{s \in \mathcal{H}} \quad \|s\|_{\mathcal{H}}^2 \tag{18a}$$

$$\text{subj. to} \quad |s(x_n) - y_n| \leq \bar{\delta}_n \tag{18b}$$

$$\forall n = 1, \ldots, N$$

where $\bar{\delta}_n$ is the $n$th element of the $\bar{\delta}$ vector. The mathematical program above can be interpreted as a $\varepsilon$-support vector regression (SVR) problem with hard margins (Schölkopf et al., 2018).

With the help of an indicator function, the constraints above can be incorporated into the objective, permitting the use of the representer theorem once more (Schölkopf et al., 2001). The optimizer can be found by solving the simple quadratic program

$$\alpha^\star = \underset{\alpha \in \mathbb{R}^N}{\arg\min} \quad \alpha^\top K \alpha \tag{19a}$$

$$\text{subj. to} \quad |K\alpha - y| \leq \bar{\delta} \tag{19b}$$

and setting $s^\star(x) = \alpha^{\star\top} K_{Xx}$, where the attained values at the data sites are denoted by $d_n := s^\star(x_n)$, $n = 1, \ldots, N$, with $d = \begin{bmatrix} d_1 & \ldots & d_N \end{bmatrix}^\top = K\alpha^\star$. Since the ground-truth is a feasible solution to (19), the order $\|s^\star\|_\mathcal{H} \leq \|f\|_\mathcal{H} \leq \Gamma$ holds. Furthermore, we highlight that

$$\|s^\star\|_\mathcal{H} \leq \|\bar{s}\|_\mathcal{H} \tag{20}$$

which also follows from the fact that the function $\bar{s}$, interpolating the noise free samples, is contained in the search space of (19). This enables us to bound the prediction-error associated with the SVR model $s^\star$.

**Theorem 2** *Let $\bar{\delta} \in \mathbb{R}_{>0}^N$ be the noise bound. Under Assumptions 1 to 3, the SVR model $s^\star$ admits the error bound*

$$|s^\star(x) - f(x)| \leq P(x)\sqrt{\Gamma^2 - \|s^\star\|_\mathcal{H}^2} + \bar{\delta}^\top |K^{-1} K_{Xx}| + |(d - y)^\top K^{-1} K_{Xx}| \tag{21}$$

*for all $x \in \mathcal{X}$, where $f$ is the unknown ground-truth and $\|s^\star\|_\mathcal{H}^2 = d^\top K^{-1} d$.*

**Proof** Observe that $d - f_X = d - y + \delta$, with $|\delta| \leq \bar{\delta}$ by Assumption 2 and where $d - y$ is known. We get

$$|s^\star(x) - f(x)| = |d^\top K^{-1} K_{Xx} - f(x)| \tag{22}$$

$$\leq |f_X^\top K^{-1} K_{Xx} - f(x)| + |(d - f_X)^\top K^{-1} K_{Xx}| \tag{23}$$

$$\leq P(x)\sqrt{\Gamma^2 - \|\bar{s}\|_\mathcal{H}^2} + |(d - y + \delta)^\top K^{-1} K_{Xx}| \tag{24}$$

$$\leq P(x)\sqrt{\Gamma^2 - \|s^\star\|_\mathcal{H}^2} + |(d - y)^\top K^{-1} K_{Xx}| + \bar{\delta}^\top |K^{-1} K_{Xx}| \tag{25}$$

Where (23) follows from the triangle inequality, (24) from Proposition 1 and the observation, and (25) from (20), the triangle inequality and the noise bound. ∎

We notice again that this bound is consistent with the noise-free case: for $\bar{\delta} = \mathbf{0}$, it holds that $d_n = f(x_n)$ and $\|s^\star\|_\mathcal{H} = \|\bar{s}\|_\mathcal{H}$, so we recover the bound in Proposition 1.

## 4. Discussion

### 4.1. Comparing KRR and SVR

Whether a KRR or SVR model should be used in a given context depends on a number of practical considerations. Strengths and weaknesses of each approach are examined next.

*Model computation:* The KRR model can be directly constructed from a dataset $D$ as in (13), since (12) admits a closed-form solution. The SVR problem (18), on the other hand, does not have an explicit solution in terms of the data, but it requires the user to solve (19).

*Hyperparameters:* The choice of the regularization parameter $\lambda$ in KRR is not straightforward, and should be guided by the knowledge one has on $\delta$. A badly chosen regularizer impacts not only the model quality, but also the size of the error bounds (15). In the Gaussian processes field, hyperparameter learning (including the kernel constants) is typically done through maximal likelihood estimation (Williams and Rasmussen, 2006); in the non-Bayesian kernel literature, different forms of cross validation are usually employed (Duan et al., 2003). In contrast with KRR, no hyperparameter is involved in the SVR alternative, making it more suitable to scenarios where nothing is known about the possible noise realization.

*Bound computation:* To compute $\Delta$ in (15), an optimization problem has to be solved. Since $\Delta$ does not depend on the query point location $x$, this has to be done only once for a fixed dataset $D$ and noise bound $\bar{\delta}$. As opposed to it, bounding the prediction error of an SVR model can be done directly as the inequality (21) only depends on given, fixed quantities.

*Error bounds magnitude:* Differences in magnitude and shape between the error bounds (15) and (21) are mainly due to their last absolute value terms. They are related to the difference of the respective kernel model and the interpolant $\tilde{s}$. In principle, choosing $\lambda$ small in the KRR setting would lead to tighter bounds. Nevertheless, this would also deteriorate the nominal model $s^*$ performance as it would fit noise rather than filtering it. An appropriate $\lambda$ is the key to attain a smooth nominal predictor while keeping the bounds small. The SVR is constrained to stay close to $\tilde{s}$ at the sample locations by definition, maximizing smoothness within the available margins. Although in our experiments they performed similarly and none of the bounds strictly encompassed the other, the KRR one seems to be slightly tighter on average for a well chosen regularization weight $\lambda$.

### 4.2. The effects of incorporating new data

The problem of improving the reconstruction quality of a surrogate model is central in approximation theory (Iske, 2000; De Marchi, 2003; De Marchi et al., 2005). Three main tools are usually employed to analyze it: the Lebesgue function $\mathcal{L}(x) := \|K^{-1}K_{Xx}\|_1$, the fill distance and the separation distance, respectively defined as

$$h_{D,\mathcal{X}} := \sup_{x \in \mathcal{X}} \min_{x_n \in D} \|x - x_n\|_2$$

$$q_D := \min_{\substack{x_i, x_j \in D \\ x_i \neq x_j}} \frac{1}{2} \|x_i - x_j\|_2$$

Notice that, if the noise bound is uniform across all samples, then $\bar{\delta}^\top |K^{-1}K_{Xx}|$ present both in (15) and (21) simplifies to $\mathcal{L}(x)$ times the bounding constant. Due to this term, it is not guaranteed that the bounds will shrink everywhere after the addition of new data-points at arbitrary locations. It is known that a new datum is most benign when it minimizes $h_{D,\mathcal{X}}$ while *not reducing* $q_D$. Balancing

these two constants is an issue commonly referred to as the *uncertainty principle* Fasshauer (2011). A key and simple advice is to use a uniformly or quasi-uniformly distributed dataset, which not only favors the bound shrinkage, but also controls the increase of the kernel matrix $K$ condition number Hangelbroek et al. (2010); Diederichs and Iske (2019). This suggests that if the data at hand are highly scattered, a pre-processing stage is highly recommended, possibly dropping points that are too close to each other as they could lead to numerical instabilities.

Despite the statements made above, and using the same arguments as in (Hashimoto et al., 2020, Sec. 4.1), one sees that the space to where the unknown ground-truth is confined is non-increasing with the addition of any new datum. In more contrete terms, denote the right-hand side of (15) by $e(x)$, and let $W_D(x) := \{y \in \mathbb{R} \,|\, s^*(x) - e(x) \leq y \leq s^*(x) + e(x)\}$ be the interval function that bounds the value of $f(x)$ for all $x \in \mathcal{X}$. Let $D^+$ be the dataset augmented with one (pairwise-distinct) point. Whereas it is not guaranteed that $W_{D^+}(x) \subset W_D(x)$, clearly $f$ must satisfy $f(x) \in W_D(x) \cap W_{D^+}(x)$, so that the confinement space is always non-increasing, i.e., total information gathered about the unknown function increased. The same observations are clearly true for the SVR model as well.

### 4.3. On the connections between KRR and Gaussian processes

Consider without loss of generality a Gaussian process setting with a null prior mean function. Let $\sigma_\delta^2$ denote the variance of the Gaussian measurement noise. Then, the GP conditional expectation defines *exactly the same map* from $\mathcal{X} \to \mathbb{R}$ as the KRR solution $s^*$ with a regularization parameter $\lambda = \sigma_\delta^2/N$ (Kanagawa et al., 2018, Sec. 3). Indeed, GP posterior means have precisely the same form as (13). The disparity between both methodologies is in the way their hypotheses spaces are defined. On one hand, $f$ is assumed to follow a distribution governed by the covariance kernel $k$, and on the other, $f$ is a static member of the RKHS of the kernel $k$. In other words, in the former case $f$ is a stochastic process realization in a suitable probability space, and in the latter, $f$ is a fixed map. Bounds derived for each model are therefore very different in nature: GP results draw probabilistic limits for their sample paths, whereas KRR results bound any function in $\mathcal{H}$, including the ground-truth Kanagawa et al. (2018).

In adopting the perspective presented in this paper, one might be concerned with the restrictiveness of working only in the space $\mathcal{H}$ (Lederer et al., 2019). It has been shown that, whereas a GP mean belongs to its reproducing kernel Hilbert space, sample paths fall outside of it almost surely as they are 'more complex' maps (Kanagawa et al., 2018, Sec. 4). Nonetheless, given any continuous function $g$, the set $\mathcal{H}$ associated with for instance the squared-exponential kernel has at least one member that is arbitrarily close to $g$—that is, $\mathcal{H}$ is dense in the class of continuous functions (Micchelli et al., 2006). Models $s$ in $\mathcal{H}$ enjoy therefore the so called universal approximation property, which renders their representation capabilities equal to many families of widely used neural networks.

## 5. Numerical experiments

Two examples are presented here to illustrate the behavior of the established inequalities in different conditions[2]. First we compare the KRR and SVR approaches to each other, and to the deterministic

---

2. The code to reproduce the examples is available at `https://github.com/emilioMaddalena/DetErrBnd`.
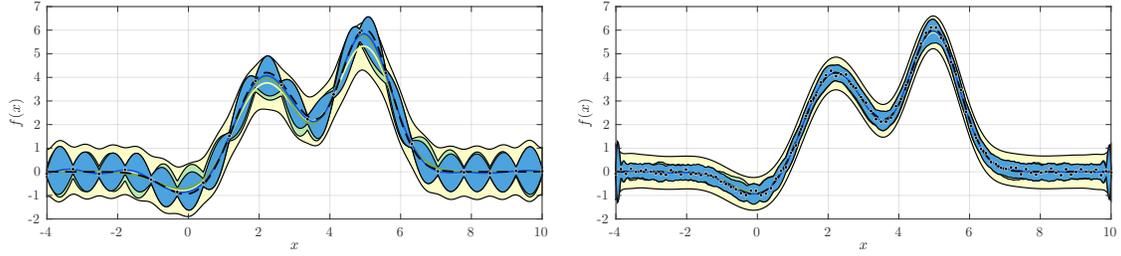
Figure 1: Ground-truth (- -), KRR (—), SVR (—), and the alternative proposed in Hashimoto et al. (2020) (—). The error bounds are depicted using the same colors of their respective models, and were computed for $N = 20$ (top) and $N = 100$ (bottom) samples. The noisy data-points are shown as black circles.

error bound recently proposed in Hashimoto et al. (2020)[3]. In the cited work, the authors proposed a GP-like surrogate model and established deterministic error bounds also exploiting an RKHS norm estimate $\Gamma$. The expression is reproduced here using our notation:

$$|s^\diamond(x) - f(x)| \le \sigma(x) \sqrt{\Gamma^2 - y^\top (K + \tilde{\delta}^2 I)^{-1} y + N} \tag{26}$$

where $s^\diamond$ is their model, $\sigma^2(x) = k(x,x) - K_{xX}(K + \tilde{\delta}^2 I)^{-1} K_{Xx}$, and $\tilde{\delta}$ is a *necessarily uniform* bound on the noise. A second example is also discussed, focusing exclusively on the KRR case. The negative effects of having highly scattered data are highlighted and a simple way to handle them is shown.

## 5.1. A comparison among three alternatives

Let $k$ be the isotropic squared-exponential kernel function

$$k(x, x_n) = \exp\left(-\frac{\|x - x_n\|_2^2}{2\ell^2}\right) \tag{27}$$

with lengthscale $\ell = 0.707$, and consider a ground-truth function $f : \mathcal{X} \to \mathbb{R}$ with a known kernel expansion $f(x) = -k(x,0) + 3.5\,k(x,2) + 1.6\,k(x,3) + 6\,k(x,5)$, which leads us to $\|f\|_{\mathcal{H}} = 7.49$. Surrogate models were built based on an overestimate $\Gamma = 9$, and two datasets sampled from $\mathcal{X} = \{x \in \mathbb{R}\,|\,-4 \le x \le 10\}$, affected by a uniformly bounded noise $|\delta_n| \le 0.15$, for all $n$. The ridge regression (with $\lambda = 0.001$), support vector regression and the model in Hashimoto et al. (2020) were computed and are shown in Figure 1. In terms of nominal predictions, the three approaches yielded similar results. The SVR model however was able to filter the existing noise best (for this specific KRR choice of $\lambda$). As for the bounds, the technique proposed in Hashimoto et al. (2020) encompassed the KRR and SVR areas almost everywhere in both scenarios, therefore being more conservative. In this particular example, this was due to two reasons. Firstly, recall Definition 2, and notice that the power function $P(x)$ will always evaluate to a number smaller than $\sigma(x)$ in (26) due to $(K + \tilde{\delta}^2 I) \succ K$. Secondly, (26) has a direct dependence on the number of

---

3. The results were derived from the well-known paper Srinivas et al. (2012).
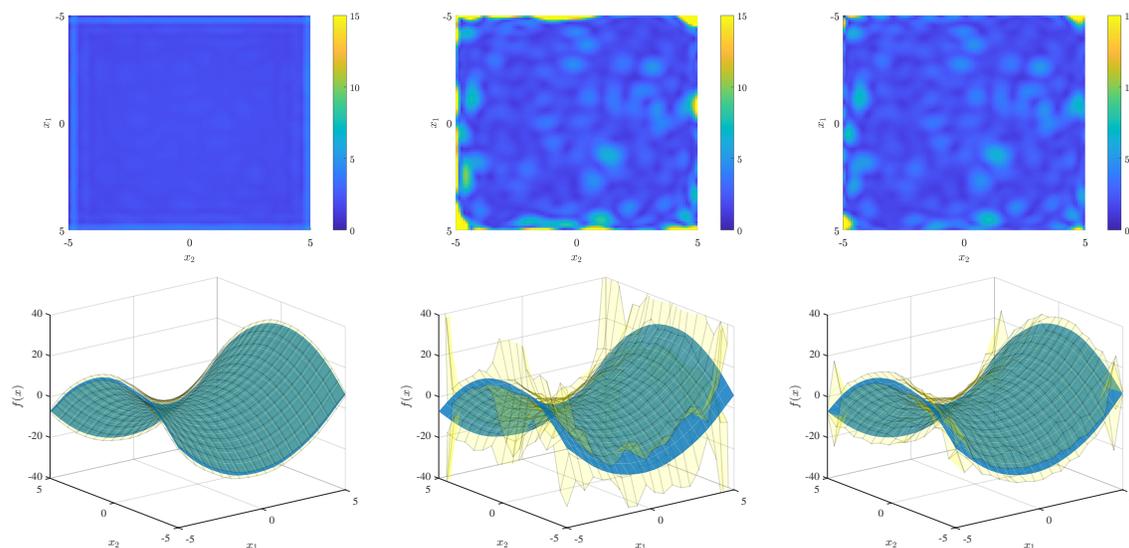
Figure 2: Error bound magnitude $e(x)$ (colormap, top) and KRR regressors along with their prediction envelopes (bottom) for three scenarios: (left) $N = 625$ points distributed in a grid, (center) $N = 625$ sampled from a uniform random distribution, and (right) $N = 625$ sampled from a uniform random distribution plus 44 points collected from the edges.

samples $N$, which is not present in the other two bounds. The uncertainty regions defined by the KRR and SVR models were similar to a large extent in shape and magnitude, indicating that both are equally suitable modeling tools. When the number of samples was increased from $N = 20$ to $N = 100$, all bounds tended to shrink: the average thickness of the prediction envelope was reduced from 2.14 to 1.41 (yellow), 1.20 to 0.73 (blue), and 1.35 to 0.74 (green). However, one can see a growth of the blue and green bounds at the extremes of the domain, when $N$ was increased. This is a common effect and was mainly caused by the Lebesgue-related term present in (15) and (21) (see Subsection 4.2).

## 5.2. The KRR bounds in 2 dimensions

Consider the dynamics of the Tinkerbell chaotic system first coordinate $f(x) = x_1^2 - x_2^2 + 0.8\,x_1 - 0.6\,x_2$ on the domain $\mathcal{X} = \{x \in \mathbb{R}^2 \mid \begin{bmatrix} -5 & -5 \end{bmatrix}^\top \le x \le \begin{bmatrix} 5 & 5 \end{bmatrix}^\top\}$. With a slight abuse of notation, $x_1$ and $x_2$ denote the first and second components of $x$. Two training datasets of $N = 625$ points were collected: one forming a perfect grid across the domain, and one drawn randomly from a uniform distribution. Bounded measurement noise $|\delta_n| \le 0.5$ for all $n$ was considered in both cases. The squared-exponential kernel was chosen and the lengthscale $\ell = 1.62$ was determined by maximizing the resulting log-likelihood for a sensible variance value. $\Gamma$ was estimated by collecting noiseless evaluations $f_X$ of the ground-truth and determining the norm of the associated interpolant; a final value of $\Gamma = 196.1$ was adopted after the use of a safety factor (see Remark 2). The KRR model (13) with $\lambda = 1 \times 10^{-5}$ was used to reconstruct $f$. The final surrogate functions $s^*(x)$ and their error bounds $e(x)$ (the right-hand side of inequality (15)) are shown in Figure 2.

As can be seen from the first two surfaces, the nominal KRR models were very similar despite being trained with differently distributed datasets. The bounds were tight and uniform under the evenly spaced samples, whereas they behaved badly under the scattered ones, showing high peaks especially at the border of the domain. We stress that the scale was held constant (from $0$ to $15$) across $e(x)$ plots for visualization purposes, but the attained values were higher in the completely saturated yellow regions as indicated by the prediction envelope below it. Notice nevertheless that the center part of the error bounds remained relatively tight. Finally, the random dataset was augmented with $44$ points collected from the domain boundary, and the results are presented in the rightmost plots. Incorporating these extra points was enough to significantly dampen the bounds increase not only at the borders, but also in internal regions. We observed an average error bound of $2.10$ for the grid sampling, an increase to $3.60$ after randomizing the sample locations, and a reduction to $2.66$ in the final case, providing thus a clear overall improvement.

## 6. Final remarks

Deterministic prediction-error bounds were provided for two classes of popular non-parametric kernel techniques: kernel ridge regression and $\varepsilon$-support vector regression. In our setting, we considered bounded measurement noise and assumed that the ground-truth function belongs to the RKHS induced by a known positive-definite kernel. We believe that our uncertainty bounds can be employed in a number of different scenarios such as in the deterministic certification of machine learning algorithms. The inequalities established herein are centered around each of the models; a pertinent question is whether the maps $f \in \mathcal{H}$, $\|f\|_{\mathcal{H}} \leq \Gamma$ alone also admit (non-trivial) bounds for their point-wise evaluations. As a last observation, the matter of hyperparameter selection is still open: what specific metrics could be considered when learning them in our bounded noise setting?

## Appendix A. Proof of Proposition 2

In the following, we only prove $f_X^{\max}(x) \geq f_Z^{\max}(x)$, the inequality for $f^{\min}$ follows from the same arguments, which proves the interval containment. To get the interval at a point $x \in \mathcal{X}$, we consider the sets $X, \bar{X} = X \cup \{x\}, Z = X \cup \{z\}$ and $W = Z \cup \{x\}$. Additionally, we denote the interpolant of $f_X$ by $s_X$ and follow this convention for the other sets. We observe the following norm identities, derived as in the proof of Proposition 1

$$\|s_W\|_{\mathcal{H}}^2 = \|s_{\bar{X}}\|_{\mathcal{H}}^2 + P_{\bar{X}}^{-2}(z)(s_{\bar{X}}(z) - f_z)^2 \tag{28}$$

$$= \|s_X\|_{\mathcal{H}}^2 + P_X^{-2}(x)(s_X(x) - f_x)^2 + P_{\bar{X}}^{-2}(z)(s_{\bar{X}}(z) - f_z)^2 \tag{29}$$

$$= \|s_Z\|_{\mathcal{H}}^2 + P_Z^{-2}(x)(s_Z(x) - f_x)^2 \tag{30}$$

$$\leq \Gamma^2 \tag{31}$$

This allows us to write $f_Z^{\max}(x)$ in two different ways

$$f_Z^{\max}(x) = s_Z(x) + P_Z(x)\sqrt{\Gamma^2 - \|s_Z\|_{\mathcal{H}}^2} \tag{32}$$

$$= s_X(x) + P_X(x)\sqrt{\Gamma^2 - \|s_X\|_{\mathcal{H}}^2 - P_{\bar{X}}^2(s_{\bar{X}}(z) - f_z)^2} \tag{33}$$

From (33), we observe

$$f_Z^{\max}(x) \leq s_X(x) + P_X(x)\sqrt{\Gamma^2 - \|s_X\|_{\mathcal{H}}^2} = f_X^{\max}(x) \tag{34}$$

using the positivity of the power function. □

## Acknowledgments

## References

N. Aronszajn. Theory of reproducing kernels. *Transactions of the American mathematical society*, 68(3):337–404, 1950.

M. Binder, G. Darivianakis, A. Eichler, and J. Lygeros. Approximate explicit model predictive controller using Gaussian processes. In *IEEE Conference on Decision and Control (CDC)*, pages 841–846. IEEE, 2019.

L. Blanken and T. Oomen. Kernel-based identification of non-causal systems with application to inverse model control. *Automatica*, 114:108830, 2020.

E. Bradford, L. Imsland, and E. A. del Rio-Chanona. Nonlinear model predictive control with explicit back-offs for Gaussian process state space models. In *IEEE Conference on Decision and Control (CDC)*, pages 4747–4754. IEEE, 2019.

E. Bradford, L. Imsland, D. Zhang, and E. A. del Rio Chanona. Stochastic data-driven model predictive control using Gaussian processes. *Computers & Chemical Engineering*, 139:1–18, 2020.

F. P. Carli, T. Chen, and L. Ljung. Maximum entropy kernels for system identification. *IEEE Transactions on Automatic Control*, 62(3):1471–1477, 2016.

T. Chen. Continuous-time DC kernel—a stable generalized first-order spline kernel. *IEEE Transactions on Automatic Control*, 63(12):4442–4447, 2018.

A. Chiuso and G. Pillonetto. System identification: A machine learning perspective. *Annual Review of Control, Robotics, and Autonomous Systems*, 2:281–304, 2019.

S De Marchi. On optimal center locations for radial basis function interpolation: computational aspects. *Rend. Splines Radial Basis Functions and Applications*, 61(3):343–358, 2003.

S. De Marchi, R. Schaback, and H. Wendland. Near-optimal data-independent point locations for radial basis function interpolation. *Advances in Computational Mathematics*, 23(3):317–330, 2005.

B. Diederichs and A. Iske. Improved estimates for condition numbers of radial basis function interpolation matrices. *Journal of Approximation Theory*, 238:38–51, 2019.

S. Diwale and C. Jones. A generalized representer theorem for hilbert space-valued functions. *arXiv preprint arXiv:1809.07347*, 2018.

A. d'Onofrio. *Bounded noises in physics, biology, and engineering*. Springer, 2013.

K. Duan, S. S. Keerthi, and A. N. Poo. Evaluation of simple performance measures for tuning SVM hyperparameters. *Neurocomputing*, 51:41–59, 2003.

G. E. Fasshauer. Positive definite kernels: past, present and future. *Dolomites Research Notes on Approximation*, 4:21–63, 2011.

T. Hangelbroek, F. J. Narcowich, and J. D. Ward. Kernel approximation on manifolds i: bounding the lebesgue constant. *SIAM Journal on Mathematical Analysis*, 42(4):1732–1760, 2010.

K. Hashimoto, A. Saoud, M. Kishida, T. Ushio, and D. Dimarogonas. Learning-based safe symbolic abstractions for nonlinear control systems. *arXiv preprint arXiv:2004.01879*, 2020.

L. Hewing, K. P. Wabersich, M. Menner, and M. N. Zeilinger. Learning-based model predictive control: Toward safe learning in control. *Annual Review of Control, Robotics, and Autonomous Systems*, 3:269–296, 2020.

A. Iske. Optimal distribution of centers for radial basis function methods. 2000.

J. Jackson, L. Laurenti, E. Frew, and M. Lahijanian. Safety verification of unknown dynamical systems via gaussian process regression. *arXiv preprint arXiv:2004.01821*, 2020.

V. Jakhetiya, K. Gu, S. Jaiswal, T. Singhal, and Z. Xia. Kernel ridge regression based quality measure and enhancement of 3D-synthesized images. *IEEE Transactions on Industrial Electronics*, 68(1):423–433, 2020.

M. Kanagawa, P. Hennig, D. Sejdinovic, and B. K. Sriperumbudur. Gaussian processes and kernel methods: A review on connections and equivalences. *arXiv preprint arXiv:1807.02582*, 2018.

T. Koller, F. Berkenkamp, M. Turchetta, and A. Krause. Learning-based model predictive control for safe exploration. In *IEEE Conference on Decision and Control (CDC)*, pages 6059–6066. IEEE, 2018.

J. Lataire and T. Chen. Transfer function and transient estimation by Gaussian process regression in the frequency domain. *Automatica*, 72:217–229, 2016.

A. Lederer, J. Umlauft, and S. Hirche. Uniform error bounds for gaussian process regression with application to safe control. In *Advances in Neural Information Processing Systems*, pages 659–669, 2019.

L. Ljung, T. Chen, and B. Mu. A shift in paradigm for system identification. *International Journal of Control*, 93(2):173–180, 2020.

J. H. Manton and P.-O. Amblard. A primer on reproducing kernel Hilbert spaces. *Foundations and Trends® in Signal Processing*, 8(1–2):1–126, 2015.

J. M. Manzano, D. Limon, D. de la Peña, and J.-P. Calliess. Robust learning-based MPC for nonlinear constrained systems. *Automatica*, 117:108948, 2020.

C. A. Micchelli, Y. Xu, and H. Zhang. Universal kernels. *Journal of Machine Learning Research*, 7(Dec):2651–2667, 2006.

R. Moriconi, K. S. S. Kumar, and M. P. Deisenroth. High-dimensional Bayesian optimization with projections using quantile Gaussian processes. *Optimization Letters*, 14(1):51–64, 2020.

C. Novara, A. Nicolì, and G. C. Calafiore. Nonlinear system identification in sobolev spaces. *arXiv preprint arXiv:1911.02930*, 2019.

G. Pillonetto and G. De Nicolao. A new kernel-based approach for linear system identification. *Automatica*, 46(1):81–93, 2010.

G. Pillonetto, M. H. Quang, and A. Chiuso. A new kernel-based approach for nonlinear system identification. *IEEE Transactions on Automatic Control*, 56(12):2825–2840, 2011.

G. Pillonetto, F. Dinuzzo, T. Chen, G. De Nicolao, and L. Ljung. Kernel methods in system identification, machine learning and function estimation: A survey. *Automatica*, 50(3):657–682, 2014.

K. Polymenakos, L. Laurenti, A. Patane, J.-P. Calliess, L. Cardelli, M. Kwiatkowska, A. Abate, and S. Roberts. Safety guarantees for planning based on iterative gaussian processes. *arXiv preprint arXiv:1912.00071*, 2019.

R. S. Risuleo, G. Bottegal, and H. Hjalmarsson. A nonparametric kernel-based approach to Hammerstein system identification. *Automatica*, 85:234–247, 2017.

R. S. Risuleo, F. Lindsten, and H. Hjalmarsson. Bayesian nonparametric identification of wiener systems. *Automatica*, 108:1–8, 2019.

U. Rosolia, X. Zhang, and F. Borrelli. Robust learning model predictive control for iterative tasks: Learning from experience. In *IEEE Conference on Decision and Control (CDC)*, pages 1157–1162, 2017.

M. Scandella, M. Mazzoleni, S. Formentin, and F. Previdi. A note on the numerical solutions of kernel-based learning problems. *IEEE Transactions on Automatic Control*, 2020.

R. Schaback. A unified theory of radial basis functions: Native hilbert spaces for radial basis functions ii. *Journal of computational and applied mathematics*, 121(1-2):165–177, 2000.

B. Schölkopf, R. Herbrich, and A. J. Smola. A generalized representer theorem. In *International conference on computational learning theory*, pages 416–426. Springer, 2001.

B. Schölkopf, A. J. Smola, and F. Bach. *Learning with kernels: support vector machines, regularization, optimization, and beyond*. the MIT Press, 2018.

U. K. Singh, R. Mitra, V. Bhatia, and A. K. Mishra. Kernel LMS-based estimation techniques for radar systems. *IEEE Transactions on Aerospace and Electronic Systems*, 55(5):2501–2515, 2019.

N. Srinivas, A. Krause, S. M. Kakade, and M. W. Seeger. Information-theoretic regret bounds for Gaussian process optimization in the bandit setting. *IEEE Transactions on Information Theory*, 58(5):3250–3265, 2012.

J. Umlauft and S. Hirche. Feedback linearization based on Gaussian processes with event-triggered online learning. *IEEE Transactions on Automatic Control*, 65(10):4154–4169, 2019.

J. Umlauft and S. Hirche. Learning stochastically stable Gaussian process state-space models. *IFAC Journal of Systems and Control*, 12:1–15, 2020.

W. Wang, R. Tuo, and C. F. Jeff Wu. On prediction properties of kriging: Uniform error bounds and robustness. *Journal of the American Statistical Association*, pages 1–27, 2019.

H. Wendland. *Scattered data approximation*. Cambridge Monographs on Applied and Computational Mathematics. Cambridge University Press, 2004. ISBN 9781139456654.

Christopher KI Williams and Carl Edward Rasmussen. *Gaussian processes for machine learning*, volume 2. MIT press Cambridge, MA, 2006.