

Robust Uncertainty Bounds in Reproducing Kernel Hilbert Spaces: A Convex Optimization Approach

Paul Scharnhorst*, Emilio T. Maddalena*, Yuning Jiang, Colin N. Jones

Abstract—The problem of establishing out-of-sample bounds for the values of an unknown ground-truth function is considered. Kernels and their associated Hilbert spaces are the main formalism employed herein along with an observational model where outputs are corrupted by bounded measurement noise. The noise can originate from any compactly supported distribution and no independence assumptions are made on the available data. In this setting, we show how computing tight, finite-sample uncertainty bounds amounts to solving parametric quadratically constrained linear programs. Next, properties of our approach are established and its relationship with another methods is studied. Numerical experiments are presented to exemplify how the theory can be applied in a number of scenarios, and to contrast it with other closed-form alternatives.

Index Terms—Uncertainty bounds, reproducing kernel Hilbert space, robust guarantees.

I. INTRODUCTION

We consider the problem of quantifying the uncertainty associated with point-evaluations of an unknown ground-truth map given a dataset of observations and assumptions on its nature. The analysis differs from widespread concentration bounds found in the machine learning literature as no assumptions are made on the statistical independence of samples. This agnosticism is central when dealing with systems that incorporate memory such as physical plants that evolve in a dynamical and hence strongly correlated fashion—see [1] for a thorough discussion about situations where the typical i.i.d. premise is inadequate. In exchange, we pose conditions on the ground-truth, requiring it to belong to a specific class of functions [2], and allow for observations to be scattered, not necessarily being drawn from any specific distribution. This is rather customary in the established field of approximation theory [3]–[5]. See also [6] for a recent perspective on the advantages offered by approximation-type bounds.

The setting in this paper is that of kernel learning, which is among the most prominent modern frameworks for both classification and regression. These non-parametric techniques are usually more data-efficient than deep network architectures, and recently intriguing connections between these two methodologies were established [7]–[9]. Kernels can be regarded as similarity measures between examples in a certain feature space [10]. This space is known as the reproducing kernel Hilbert space (RKHS) and is usually an infinite-dimensional linear space of functions. Moreover, it is well

known that the RKHS associated with certain kernel classes is dense in the space of continuous functions on compact domains [11]. By requiring the latent ground-truth to be a member of the RKHS associated with a known kernel, straightforward error-bounds can be established for models that interpolate noise-free data-points (see for instance [12]). Recently, these out-of-sample guarantees were extended to regularized smoothing models in the presence of bounded measurement noise [13]. Nevertheless, the task of *exactly* quantifying the associated uncertainty in the latter scenario remained open, as well as understanding how much conservatism is introduced when centering the bounds around pre-specified models.

Contributions: Herein we investigate the uncertainty quantification problem in RKHSs and with datasets corrupted by measurement noise. The sources of uncertainty are both epistemic and aleatoric [14] as explained next. The first stems from the ground-truth being an unknown fixed member of our function class, and from which we derive information indirectly through its samples. Secondly, the additive bounded measurement noise, which could originate from any probability measure, or even be a constant, fixed bias. In contrast with the study in [13], we carry out an *algorithmic independent* analysis that is not centered around any specific model; this is done by computing the highest and lowest possible point-evaluations that are consistent with our knowledge. Our main result is to show how this infinite-dimensional problem can be translated into a finite convex quadratically constrained linear program (QCLP) without any conservatism, which is accomplished through a representer theorem. Next, properties of this procedure are derived and connections with closed-form sub-optimal bounds [13] as well as classical noise-free bounds [4] are established. Finally, efficient solution methods are proposed through the dual optimization problem, which trade-off computational time and precision. Numerical experiments are reported to illustrate their use, as well as the influence of the input distribution on the final results.

Relevance for automatic control: The use of the so-called data-driven techniques to refine models, improve performance on-line, or approximate controllers is becoming evermore present in the field of automatic control [15]–[17]. In the particular case of kernel surrogate models, a considerable body of rigorous literature exists for linear dynamics [18]–[21], linear parameter-varying dynamics [22], [23], Hammerstein and Wiener cascaded systems [24], [25], mainly adopting a time-domain perspective of the identification problem. When operational constraints are present, one has to pair these tools with appropriate uncertainty quantification techniques to not make unsafe decisions. Examples include system simulation with guaranteed accuracy [26] and controller tuning algorithms

*The first two authors contributed equally. This work has received support from the Swiss National Science Foundation under the RISK project (Risk Aware Data-Driven Demand Response, grant number 200021 175627), and CSEM's Data Program. (Corresponding author: Yuning Jiang)

All authors are with Automatic Control Laboratory, EPFL, Lausanne, Switzerland. Paul Scharnhorst is with CSEM, Neuchâtel, Switzerland. (e-mail: paul.scharnhorst, emilio.maddalena, yuning.jiang, colin.jones@epfl.ch)

that avoid unreliable parameters [27], [28]. By establishing our optimal, non-asymptotic uncertainty bounds, our work aims at bridging non-parametric kernel learning and robust analysis and control. Practical applications of the results include predictive control schemes that explicitly incorporate non-parametric uncertainties [29], [30], and the certification of machine learning-based algorithms [31], [32], but in a deterministic fashion. More generally, the provided tools could also be employed in the domain of real-time optimization under unknown constraints [33]. Our motivation is similar in essence to the ones found in non-linear set membership and interval analysis works [34], [35], but our study is focused on kernel machines and their associated spaces.

Notation: \mathbb{N} denotes the set of natural numbers and \mathbb{R}^d is the d -dimensional Euclidean space. Let S_1 and S_2 be subspaces of S , then $S_1 \oplus S_2 = S$ indicates their vector direct sum, i.e., $\forall s \in S, \exists! s_1 \in S_1, \exists! s_2 \in S_2 : s = s_1 + s_2$. We denote by K_{XX} the matrix of kernel evaluations at X , i.e., the matrix containing $k(x_i, x_j)$ at its i -th row and j -th column, $x_i, x_j \in X$. Let a query point x be specified, then K_{Xx} represents the column vector-valued function $x \mapsto [k(x, x_1) \ \dots \ k(x, x_d)]^\top \in \mathbb{R}^d$, whereas K_{xX} denotes its transpose. For a matrix A we denote its nullspace by $N(A)$.

II. PRELIMINARY: KERNELS AND THEIR RKHS

We start by briefly reviewing the formalism of kernel learning, and define the main elements that will be later used in our analysis. The reader is referred to [12], [36] for further details on this topic.

A kernel $k : \Omega \times \Omega \rightarrow \mathbb{R}$ is any symmetric, real-valued function defined on a non-empty input set Ω . k is said to be positive-definite if the weighted sum $\sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j k(x_i, x_j)$ is strictly positive $\forall n \in \mathbb{N}, \forall \alpha_1, \dots, \alpha_n \in \mathbb{R} \setminus \{0\}, \forall x_1, \dots, x_n \in \Omega$. An example of a commonly used continuous kernel that enjoys this property is the squared-exponential, also known as the radial basis function (RBF) kernel. Associated with each k , there is a unique Hilbert space of maps \mathcal{H} that is referred to as the reproducing kernel Hilbert space (RKHS) of k . For compact domains Ω , [11] presents families of kernels whose \mathcal{H} are dense in the space of continuous functions, which can be interpreted as a measure of richness of such a space. Let $\mathbb{R}^\Omega = \{f : \Omega \rightarrow \mathbb{R}\}$ and $L_x : \mathbb{R}^\Omega \rightarrow \mathbb{R}$ be the map $L_x : f \mapsto f(x)$, also known as the evaluation functional for a given $x \in \Omega$. Formally, a RKHS is simply a Hilbert space $\mathcal{H} \subset \mathbb{R}^\Omega$ for which the L_x maps are continuous $\forall x \in \Omega$. It turns out that partially-evaluated kernels $k(x, \cdot)$ belong to \mathcal{H} , $\forall x \in \Omega$ and define evaluation functionals through $\langle f, k(x, \cdot) \rangle_{\mathcal{H}} = f(x)$, $\forall f \in \mathcal{H}, \forall x \in \Omega$. The latter is known as the reproducing property. From a constructive viewpoint, \mathcal{H} is given as the closure (w.r.t. the topology induced by the inner-product) of $\text{span}(\{k(x, \cdot), x \in \Omega\})$, encompassing thus weighted sums of partial kernels and limit points of sequences as well. It can be shown that this construction results in proper functions and not in equivalence classes of them.

Let $f \in \mathcal{H}$ with finite expansion $f = \sum_{i=1}^{n_f} \alpha_i k(x_i, \cdot)$, $\alpha_i \in \mathbb{R}, x_i \in \Omega$, for all i . Its induced norm $\|f\|_{\mathcal{H}}$ is then given by

$$\begin{aligned} \|f\|_{\mathcal{H}}^2 &= \left\langle \sum_{i=1}^{n_f} \alpha_i k(x_i, \cdot), \sum_{i=1}^{n_f} \alpha_i k(x_i, \cdot) \right\rangle_{\mathcal{H}} \\ &= \sum_{i=1}^{n_f} \sum_{j=1}^{n_f} \alpha_i \alpha_j k(x_i, x_j) = \alpha^\top K_{XX} \alpha \end{aligned} \quad (1)$$

due to the reproducing property, with α being the vector of scalar weights. If a member $f \in \mathcal{H}$ is the limit of a sequence $f = \sum_{i=1}^{\infty} \alpha_i k(x_i, \cdot)$, then its norm is

$$\|f\|_{\mathcal{H}}^2 = \lim_{n \rightarrow \infty} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j k(x_i, x_j).$$

As a last introductory step, we consider a finite subset $X \subset \Omega$ and define the power function $P_X : \Omega \rightarrow \mathbb{R}_{\geq 0}$ as

$$P_X(x) = \sqrt{k(x, x) - K_{xX} K_{XX}^{-1} K_{Xx}} \quad (2)$$

whenever clear from the context, the reference to X will be omitted. $P_X(x)$ can be interpreted as a form of statistical covariance, and evaluates to zero $\forall x \in X$.

III. OPTIMAL BOUNDS IN RKHS

This section introduces the main technical results of this work. First, an infinite-dimensional variational problem is formulated to bound the ground-truth values at unseen locations, and its equivalence to a finite-dimensional problem is shown. Then, we discuss properties of the derived bounds and a closed-form alternative that does not involve solving any optimization problem.

Herein we consider a positive-definite kernel $k : \Omega \times \Omega \rightarrow \mathbb{R}$ along with its corresponding RKHS $\mathcal{H} \subset \mathbb{R}^\Omega$. Our input space is taken to be a compact subset of the Euclidean space $\Omega \subset \mathbb{R}^n$. A dataset $\{(x_i, y_i)\}_{i=1}^d$ is given to us, being composed of inputs $x_i \in \Omega$ and outputs $y_i \in \mathbb{R}^{n_i}$, $y_i = [y_{i,1} \ \dots \ y_{i,n_i}]^\top$ that contain n_i scalar samples collected at the same input location x_i . The dataset carries information about an underlying ground-truth map $f^* \in \mathcal{H}$ according to

$$y_{i,j} = f^*(x_i) + \delta_{i,j} \quad (3)$$

where $\delta_{i,j}$ represents an additive measurement noise that is assumed to be uniformly bounded as stated next.

Assumption 1. *The magnitude of each noise realization $\delta_{i,j}$ is bounded by a known scalar quantity $\bar{\delta}$, i.e. $|\delta_{i,j}| \leq \bar{\delta}, \forall i, j$.*

Assumption 2. *An estimate $\Gamma \geq \|f^*\|_{\mathcal{H}}$ for the ground-truth norm is known.*

For notational convenience, we define the quantities $X := \{x_1, \dots, x_d\}$ and $y := [y_1^\top \ \dots \ y_d^\top]^\top$, which represent respectively the collection of inputs and the available outputs.

The aim is to quantify the uncertainty associated with values of the latent function f^* in the output space \mathbb{R} . We note

that, from the reproducing property and the Cauchy–Schwarz inequality, one can readily establish

$$|f(x)| = |\langle f, k(x, \cdot) \rangle_{\mathcal{H}}| \quad (4a)$$

$$\leq \|f\|_{\mathcal{H}} \|k(x, \cdot)\|_{\mathcal{H}} \quad (4b)$$

$$= \|f\|_{\mathcal{H}} \sqrt{k(x, x)} < \infty \quad (4c)$$

$\forall f \in \mathcal{H}$, including the ground-truth f^* . If the kernel k is translation-invariant, then $k(x, x)$ is constant $\forall x \in \Omega$ and (4) constitutes a uniform trivial bound for the unknown function. A reason for this inequality to be rather loose is that it does not incorporate any information provided by the outputs y or by the quantity $\bar{\delta}$, which we exploit next.

To upper bound the ground-truth values, we consider the following infinite-dimensional variational problem $\mathbb{P}0$, with the query point $x \in \Omega$ as a parameter

$$F(x) = \sup_{f \in \mathcal{H}} \{f(x) : \|f\|_{\mathcal{H}} \leq \Gamma, \|f_X - y\|_{\infty} \leq \bar{\delta}\} \quad (5)$$

where $f_X := \Lambda [f(x_1) \dots f(x_d)]^\top$ is the vector of evaluations at the input locations, which are repeated whenever multiple outputs are available at a given input. This is accomplished through Λ as defined in Appendix B. We highlight that the supremum is guaranteed to exist thanks to (4). Given a query location x , $\mathbb{P}0$ yields the tightest upper bound for $f(x)$ over all members $f \in \mathcal{H}$ of our hypothesis space that are consistent with our dataset, as well as our knowledge on the ground-truth complexity $\|f\|_{\mathcal{H}} \leq \Gamma$. Note how linking the function evaluations f_X and the outputs y plays a role analogous to conditioning stochastic processes on past observations in statistical frameworks.

Consider now the convex parametric quadratically-constrained linear program $\mathbb{P}1$

$$C(x) = \max_{c \in \mathbb{R}^d, c_x \in \mathbb{R}} c_x \quad (6a)$$

$$\text{subj. to } \begin{bmatrix} c \\ c_x \end{bmatrix}^\top \begin{bmatrix} K_{XX} & K_{Xx} \\ K_{xX} & k(x, x) \end{bmatrix}^{-1} \begin{bmatrix} c \\ c_x \end{bmatrix} \leq \Gamma^2 \quad (6b)$$

$$\|\Lambda c - y\|_{\infty} \leq \bar{\delta} \quad (6c)$$

for any $x \in \Omega \setminus \{X\}$, and extend its value function to points from the dataset $x = x_i \in X$ with the solution of $\mathbb{P}1' : C(x_i) = \max_{c \in \mathbb{R}^d} \{c_i \mid c^\top K_{XX}^{-1} c \leq \Gamma^2, \|\Lambda c - y\|_{\infty} \leq \bar{\delta}\}$, where c_i is the i -th component of c . This can be thought of as finding a map that interpolates the points $\{(x_i, c_i)\}_{i=1}^d$ and maximizes its value c_x at the input location x . The two cases $\mathbb{P}1$ and $\mathbb{P}1'$ are distinguished due to the matrix in (6b) becoming singular for any $x \in X$, and since it allows for one decision variable to be eliminated. Finally, the connection between (5) and (6) is stated next.

Theorem 1. (Finite-dimensional equivalence): *The objective in $\mathbb{P}0$ attains its supremum in \mathcal{H} and $F(x) = C(x)$ for any $x \in \Omega$.*

The derivation of the result, which is given in Appendix C-A, follows lines similar to classical representer theorem ones, i.e., showing that the optimizer necessarily lies in a finite-dimensional subspace of the RKHS. Nevertheless, note

that the objective $\mathbb{P}1$ is not regularized, nor is x necessarily an input of our dataset. Moreover, the proof also establishes the attainment property in \mathcal{H} , which helps in understanding the nature of the constraints.

Complementing (5), one could also be interested in the infimum $\inf_{f \in \mathcal{H}} \{f(x) : \|f\|_{\mathcal{H}} \leq \Gamma, \|f_X - y\|_{\infty} \leq \bar{\delta}\}$ bounding the lowest attainable value at x . In this case, a result analogous to Theorem 1 could be established, showing its equivalence to

$$B(x) = \min_{c \in \mathbb{R}^d, c_x \in \mathbb{R}} c_x \quad (7a)$$

$$\text{subj. to } \begin{bmatrix} c \\ c_x \end{bmatrix}^\top \begin{bmatrix} K_{XX} & K_{Xx} \\ K_{xX} & k(x, x) \end{bmatrix}^{-1} \begin{bmatrix} c \\ c_x \end{bmatrix} \leq \Gamma^2 \quad (7b)$$

$$\|\Lambda c - y\|_{\infty} \leq \bar{\delta} \quad (7c)$$

for any $x \in \Omega \setminus X$, and extended to $B(x_i) = \min_{c \in \mathbb{R}^d} \{c_i \mid c^\top K_{XX}^{-1} c \leq \Gamma^2, \|\Lambda c - y\|_{\infty} \leq \bar{\delta}\}$ for $x = x_i \in X$. As a result of this subsection, for any point in the domain $x \in \Omega$, the solutions to the two convex programs (6) and (7) define an *uncertainty envelope* that confines the ground-truth to its interior $B(x) \leq f^*(x) \leq C(x)$.

Remark 1. (On the necessity of Γ): Data alone are not sufficient to compute *any* out-of-sample bounds when considering functions $f \in \mathcal{H}$, regardless of the number of samples $d < \infty$ that one has. Given any tentative bound ϵ at $x \notin X$, there exists $f_\rho \in \mathcal{H}$ consistent with the dataset that will violate the bound, that is, $f_\rho(x) > \epsilon + \rho$, for any pre-specified violation level $\rho > 0$. This is simply due to the existence of maps that can interpolate any finite set of samples. Restricting the search to the Γ -ball in \mathcal{H} limits the flexibility of the considered functions, thus allowing for guarantees to be established. An analogous argument can be made in the space of Lipschitz functions. If no bound is posed on the Lipschitz constant of the ground-truth, assuming Lipschitz continuity *per se* becomes vacuous.

Remark 2. (On the noise assumption): Assumption 1 is central to the robust control literature and requires a careful handling in practice. Quantities $\bar{\delta}$ estimated from historical data could be invalidated by newly obtained samples, potentially rendering $\mathbb{P}1$ infeasible. In such a case, $\bar{\delta}$ would need to be augmented to accommodate the new samples (similar issues are discussed in [37, Section 3]). Certainly, dataset pre-processing and outlier detection are essential to the success of data-driven methods, including the present one, in practice.

Remark 3. (On having loose Γ and $\bar{\delta}$): As formulated in (5), the bound $F(x)$ depends on the quality of the available Γ and $\bar{\delta}$. The looser these two quantities are, the larger the resulting bounds—see a numerical example in Section VI. A straightforward method is presented in Appendix A to compute RKHS norm lower estimates purely based on data, which are refined the more samples one has. The methodology can help users to arrive at upper bounds Γ through augmentation.

A. Width and width shrinkage

Given our knowledge on the noise influence $\bar{\delta}$, it is natural to ask what the limits of the uncertainty quantification technique considered herein are. For example, is the width of the envelope $C(x) - B(x)$ restricted to a certain minimum value that cannot be reduced even with the addition of new data? From (6c), it is clear that at any input location $x_i \in X$, $C(x_i)$ and $B(x_i)$ cannot be more than $2\bar{\delta}$ apart. In addition to that, the presence of the complexity constraint (6b) can bring the two values closer to each other. Depending on how restrictive this latter constraint is for a given x_i , the corresponding output y_i might lie outside the interval between $C(x_i)$ and $B(x_i)$. In this case, the resulting width is considerably reduced as illustrated in Figure 1 (left).

Proposition 1. (Width smaller than the noise bound): *If $\exists y_i$ such that $y_{i,j} > C(x_i)$ or $y_{i,j} < B(x_i)$ for some j , then $C(x_i) - B(x_i) \leq \bar{\delta}$.*

Proof. Follows from $C(x_i) \geq B(x_i)$, $C(x_i) \leq y_{i,j} + \bar{\delta}$ and $B(x_i) \geq y_{i,j} - \bar{\delta}$ for any $i = 1, \dots, d$ and any $j = 1, \dots, n_i$. \square

Suppose now one has sampled (x_i, y_i) with $y_i = [y_{i,1} \ y_{i,2}]^\top$, $y_{i,1} = f^*(x_i) + \bar{\delta}$ and $y_{i,2} = f^*(x_i) - \bar{\delta}$. Then there is no uncertainty whatsoever about f^* at x_i since $f^*(x_i) = (y_{i,1} + y_{i,2})/2$ is the only possible value attainable by the ground-truth. This illustrates that the possibility of having multiple outputs at the same location allows for the uncertainty interval to shrink past the $\bar{\delta}$ width, and eventually even reduce to a singleton as shown in Figure 1 (right). Notwithstanding, the addition of a new datum to an existing dataset—be it in the form of a new output at an already sampled location or a completely new input-output pair—can only reduce the uncertainty.

Proposition 2. (Decreasing uncertainty): *Let $C_1(x)$ denote the solution of $\mathbb{P}1$ with a dataset $D_1 = \{(x_i, y_i)\}_{i=1}^d$, and $C_2(x)$ the solution with $D_2 = D_1 \cup \{(x_{d+1}, y_{d+1})\}$. Then $C_2(x) \leq C_1(x)$ for any $x \in \Omega$.*

Proof. Denote by $\mathbb{P}1_1$ the problem solved with D_1 and decision variables $[c \ c_x]$. Similarly, $\mathbb{P}1_2$ is associated with the dataset D_2 and the decision variables $[c \ c_x \ c_z]$, where c_z are due to the additional input in D_2 . Since D_2 contains all members of D_1 , the ∞ -norm constraint of $\mathbb{P}1_2$ can be recast as that of $\mathbb{P}1_1$ and an additional constraint for c_z and the new

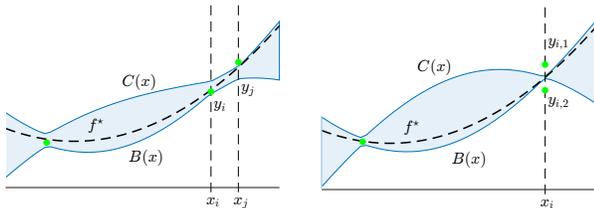


Fig. 1. (Left) A sample lying outside of the uncertainty envelope, implying that the width is smaller than $\bar{\delta}$ at x_j . (Right) Redundant information is used to shrink the uncertainty envelope. In this scenario, we recover the ground-truth value at x_i as $C(x_i) = B(x_i) = f^*(x_i)$.

outputs. Let $\mathbb{X} := X \cup \{x\}$, $\bar{c} := [c^\top \ c_x]^\top$ and $z := x_{d+1}$ be shorthand variables to ease notation. The complexity constraint of $\mathbb{P}1_2$ is then

$$\begin{bmatrix} \bar{c} \\ c_z \end{bmatrix}^\top \begin{bmatrix} K_{\mathbb{X}\mathbb{X}} & K_{\mathbb{X}z} \\ K_{z\mathbb{X}} & k(z, z) \end{bmatrix}^{-1} \begin{bmatrix} \bar{c} \\ c_z \end{bmatrix} \leq \Gamma^2 \quad (8a)$$

$$\stackrel{(i)}{\Leftrightarrow} \bar{c}^\top K_{\mathbb{X}\mathbb{X}}^{-1} \bar{c} + P_{\mathbb{X}}^{-2}(z) \left\| \begin{bmatrix} K_{\mathbb{X}\mathbb{X}}^{-1} K_{\mathbb{X}z} \\ -1 \end{bmatrix} \begin{bmatrix} \bar{c} \\ c_z \end{bmatrix} \right\|_2^2 \leq \Gamma^2 \quad (8b)$$

$$\stackrel{(ii)}{\Leftrightarrow} \begin{bmatrix} c \\ c_x \end{bmatrix}^\top \begin{bmatrix} K_{XX} & K_{Xx} \\ K_{xX} & k(x, x) \end{bmatrix}^{-1} \begin{bmatrix} c \\ c_x \end{bmatrix} + P_{\mathbb{X}}^{-2}(z) (\bar{c}^\top K_{\mathbb{X}\mathbb{X}}^{-1} K_{\mathbb{X}z} - c_z)^2 \leq \Gamma^2 \quad (8c)$$

where the matrix identity found in Appendix D was used in (i) and $P_{\mathbb{X}}^2(z) = k(z, z) - K_{z\mathbb{X}} K_{\mathbb{X}\mathbb{X}}^{-1} K_{\mathbb{X}z}$. In (ii), the definitions of \bar{c} and \mathbb{X} were used. Thanks to $P_{\mathbb{X}}(z) \geq 0, \forall z$ and the quadratic term multiplying it, we conclude that for any choice of the decision variable c_z , (8c) is a tightened version of the complexity constraint of $\mathbb{P}1_1$, which is (6b). As a result, the maximum of $\mathbb{P}1_2$ is lower or equal than that of $\mathbb{P}1_1$. \square

Let us take a closer look at the tightened constraint (8c). The term $\bar{c}^\top K_{\mathbb{X}\mathbb{X}}^{-1} K_{\mathbb{X}z} =: s(z)$ represents an interpolating model passing through the output values \bar{c} , that is, c and c_x (see e.g. the discussion in Section 3.1 of [13]). If the difference $s(z) - c_z$ can be made small, then the tightening will also be reduced, whereas it will be significant if the difference is large. The result is of course dictated by the ∞ -norm constraint, since c_z cannot be more than $\bar{\delta}$ away from all the outputs y available at z . Therefore, a new datum will cause significant shrinkage of the envelope at a point $z \in \Omega$ when the new output causes $s(z) - c_z$ to be large, which intuitively can be seen as a measure of gained information through the new sample. Finally, this process is weighted by the inverse of the power function $P_{\mathbb{X}}^{-2}(z)$, which does not depend on any output, but only on the input locations. For more practical guidelines and a visual representation of how new data can contribute to reducing the ground-truth uncertainty, the reader is referred to Examples 1 and 2 from Section VI.

Remark 4. Recovering the ground-truth as shown in Figure 1 (right) requires the noise realizations to match $\bar{\delta}$ and $-\bar{\delta}$; it is thus necessary to have tight noise bounds for it to happen. On the other hand, Proposition 2 guarantees the decreasing uncertainty property regardless of how accurate $\bar{\delta}$ is. Although not explicitly stated, a completely analogous result holds for the lower part of the envelope $B(x)$.

B. A sub-optimal closed-form alternative

The discussion in this subsection assumes that only one sample is present at each input location, i.e., $y_i = y_i$ for $i = 1, \dots, d$, so that $y = y$.

In order to alleviate the computational complexity of having to solve two optimization problems at each query point, closed-form expressions can be employed instead. These surrogates yield sub-optimal bounds around any pre-specified kernel model of the form $s(x) = \alpha^\top K_{Xx}$, for some $\alpha \in \mathbb{R}^d$.

Proposition 3. Let $s(x) = \alpha^\top K_{Xx}$, for a given $\alpha \in \mathbb{R}^d$. Then, for any $x \in \Omega$, the ground-truth is bounded by $s(x) - S(x) \leq f^*(x) \leq s(x) + S(x)$ with

$$S(x) = P_X(x) \sqrt{\Gamma^2 + \tilde{\Delta}} + \bar{\delta} \|K_{XX}^{-1} K_{Xx}\|_1 + |\tilde{s}(x) - s(x)| \quad (9)$$

where $\tilde{s}(x) = y^\top K_{XX}^{-1} K_{Xx}$, and the constant $\tilde{\Delta}$ is the minimum of the unconstrained convex problem $\min_{\nu \in \mathbb{R}^d} \{\frac{1}{4} \nu^\top K_{XX} \nu + \nu^\top y + \bar{\delta} \|\nu\|_1\}$.

Proof. See Appendix C-B.

The map $\tilde{s}(x) = y^\top K_{XX}^{-1} K_{Xx}$ is an interpolant for the available outputs y . Note also that none of the terms in (9) depend on the model weights α with the exception of the last term $|\tilde{s}(x) - s(x)|$. Therefore, the width $S(x)$ will be minimized when $s(x) = \tilde{s}(x) \implies \alpha = y^\top K_{XX}^{-1}$. Since such a model would severely overfit, a balance between smoothing the data and not diverging too much from $\tilde{s}(x)$ has to be found. In our previous work [13], we have illustrated how kernel ridge regression and minimum norm models are good candidate techniques to accomplish this goal.

By reformulating the optimal bounds, we uncover their relation with the suboptimal estimates given in Proposition 3. First, consider P1 and optimize over the decision variable $\delta = c - y$ rather than over c . Next, apply a quadratic decomposition identical to the one used in (8) to the complexity constraint (6b) and solve for c_x . After recalling that $\tilde{s}(x) = y^\top K_{XX}^{-1} K_{Xx}$ and $\|\tilde{s}\|_{\mathcal{H}}^2 = y^\top K_{XX}^{-1} y$, one obtains

$$c_x \leq P(x) \sqrt{\Gamma^2 - \|\tilde{s}\|_{\mathcal{H}}^2 - \delta^\top K_{XX}^{-1} \delta + 2y^\top K_{XX}^{-1} \delta} + \tilde{s}(x) + \delta^\top K_{XX}^{-1} K_{Xx} \quad (10)$$

Instead of maximizing c_x , the right-hand side of (10) can be directly considered as the objective function equivalently. As a result, we obtain

$$\max_{\|\delta\|_\infty \leq \bar{\delta}} \tilde{s}(x) + P(x) \sqrt{\Gamma^2 - \|\tilde{s}\|_{\mathcal{H}}^2 - \delta^\top K_{XX}^{-1} \delta - 2y^\top K_{XX}^{-1} \delta} + \delta^\top K_{XX}^{-1} K_{Xx}$$

Now, relax the problem by allowing δ to attain different values inside and outside the square-root

$$\max_{\delta_1, \delta_2 \in \mathbb{R}^d} \tilde{s}(x) + P(x) \sqrt{\Gamma^2 - \|\tilde{s}\|_{\mathcal{H}}^2 - \delta_1^\top K_{XX}^{-1} \delta_1 + 2y^\top K_{XX}^{-1} \delta_1} + \delta_2^\top K_{XX}^{-1} K_{Xx} \quad (11a)$$

$$\text{subj. to } \|\delta_1\|_\infty \leq \bar{\delta}, \|\delta_2\|_\infty \leq \bar{\delta} \quad (11b)$$

Note that the objective is separable and that $\tilde{\Delta}$ is the dual solution of $\max_{\delta_1 \in \mathbb{R}^d} \{-\delta_1^\top K_{XX}^{-1} \delta_1 + 2y^\top K_{XX}^{-1} \delta_1 - \|\tilde{s}\|_{\mathcal{H}}^2\}$. Also, $\max_{\delta_2 \in \mathbb{R}^d} \{\delta_2^\top K_{XX}^{-1} K_{Xx} : \|\delta_2\|_\infty \leq \bar{\delta}\} = \bar{\delta} \|K_{XX}^{-1} K_{Xx}\|_1$ since these norms are duals of each other. Remember that the objective (11a) is a conservative upper bound for $f^*(x)$, having $\tilde{s}(x)$ as the reference model. Given any smoother $s(x)$, the triangle inequality $|f(x) - s(x)| \leq |f(x) - \tilde{s}(x)| + |\tilde{s}(x) - s(x)|$ can be used to bound the distance between its predictions and the ground-truth values, arriving thus at the same expressions presented in Proposition 3.

From (10), the noise variable δ is seen to increase the maximum in two distinct ways: through the inner product

$\delta^\top K_{XX}^{-1} K_{Xx}$, and via a norm augmentation corresponding to $\tilde{\Delta}$. One source of conservativeness in Proposition 3 is taking into account the worst-possible inner-product and norm increase jointly. Despite this fact, they yield competitive results for moderate noise-levels as shown numerically in Section VI. We moreover note that in the noise-free scenario, (10) and (11a) are the same, and Proposition 3 simplifies to the classical bounds in the interpolation case (see for instance [38]).

Remark 5. The sub-optimal bounds presented in this subsection feature a nominal model at their center, which is desirable in many practical situations. In the optimal scenario, the minimum norm regressor $s^*(x) = \alpha^{*\top} K_{Xx}$, $\alpha^* = \arg \min_{\alpha \in \mathbb{R}^d} \{\alpha^\top K_{XX} \alpha : \|K_{XX} \alpha - y\|_\infty \leq \bar{\delta}\}$ can be used as a nominal model. This choice is guaranteed to lie completely within $C(x)$ and $B(x)$ —although not necessarily in the middle—since the map s^* belongs to \mathcal{H} and is a feasible solution for P0.

IV. EFFICIENT COMPUTATION AND OUTER APPROXIMATIONS

One of the fundamental sources of computational complexity in kernel learning lies in the inverse term K_{XX}^{-1} . Scaling these techniques to large datasets in a principled manner is still a topic of active research [39], [40]. Notice that K_{XX}^{-1} is explicitly present in P1, thus limiting its applicability to small and medium-sized problems due to the cubic time complexity associated with the inverse operation. In this section we discuss alternative formulations that can be solved more efficiently.

A. The dual problem

Following a standard dualization procedure, which can be found in Appendix C-C, a Lagrangian dual for P1 in (6) can be the convex problem D1

$$\min_{\nu \in \mathbb{R}^d, \lambda > 0} \frac{1}{4\lambda} \nu^\top \Lambda K_{XX} \Lambda^\top \nu + \left(y - \frac{1}{2\lambda} \Lambda K_{Xx} \right)^\top \nu + \bar{\delta} \|\nu\|_1 + \frac{1}{4\lambda} k(x, x) + \lambda \Gamma^2 \quad (12)$$

for any query point $x \in \Omega \setminus X$. In our notation the dimension $\tilde{d} = \sum_{i=1}^d n_i$ is the total number of outputs, that is, the size of y . As detailed in Appendix C-C, under the assumption of the complexity constraint not being active, the dual of P1' is also D1, meaning that the formulation (12) could be used $\forall x \in \Omega$. Remarkably, D1 only involves the kernel matrix itself and not its inverse, avoiding thus the aforementioned adversity. Furthermore, the query point x enters D1 through the terms K_{Xx} and $k(x, x)$. The former measures the similarity between the query point x and each of the inputs in X ; the latter is simply a constant term for translation-invariant kernels, and evaluates always to 1 in the specific case of the squared-exponential kernel.

The optimization problem above is convex since it is a quadratic-over-linear function with $\Lambda K_{XX} \Lambda^\top \succeq 0$ and λ restricted to the positive reals. The objective can moreover be decomposed into a differentiable part and a single non-differentiable term $\|\nu\|_1$, with ν unconstrained. This class

of problems has long been studied and mature numerical algorithms exist to solve them, notably different flavors of splitting methods such as the alternating direction method of multipliers (ADMM) [41, Section 6]. Alternatively, a standard linear reformulation could be employed to substitute $\|\nu\|_1$ by $\sum_i \eta_i$, with additional constraints $-\nu \leq \eta$, $\nu \leq \eta$. The result is a completely differentiable objective, but with extra decision variables and linear constraints. Next, a mild condition is given ensuring a zero duality gap between the primal and dual problems.

Proposition 4. (Strong duality): *If $\bar{\delta} > \delta_{i,j}, \forall i, j$ and $\Gamma > \|f^*\|_{\mathcal{H}}$, then no duality gap exists, i.e., $\max \mathbb{P}1 = \min \mathbb{D}1$.*

Proof. Consider the primal problem $\mathbb{P}1$ and select $c = f_X^*$ and $c_x = f^*(x)$. Let $\mathbb{X} := X \cup \{x\}$ and $K_{\mathbb{X}\mathbb{X}}$ denote the kernel matrix associated with \mathbb{X} . Thanks to the optimal recovery property [4, Theorem 13.2], $[c^\top \ c_x] K_{\mathbb{X}\mathbb{X}} [c^\top \ c_x]^\top \leq \|f^*\|_{\mathcal{H}}^2$, which in turn is strictly smaller than Γ^2 by assumption. Also, $\|\Lambda c - y\|_\infty = \|\Lambda f_X^* - y\|_\infty = \left\| \begin{bmatrix} \delta_{1,1} & \dots & \delta_{2,1} & \dots \end{bmatrix}^\top \right\|_\infty < \bar{\delta}$. Therefore, the ground-truth values constitute a feasible solution that lies in the interior of the primal problem feasible set. As a result, Slater's condition is met and, since the primal is convex, there is no duality gap. \square

B. An alternating optimization procedure

Solving the dual problem to any accuracy leads to an over bound on $C(x)$ thanks to duality. In other words, any feasible sub-optimal solution of $\mathbb{D}1$ establishes a conservative uncertainty estimate. This motivates the study of light methods that could trade-off computational time and accuracy. In what follows we describe a block coordinate minimization scheme to tackle the problem, which is later shown to yield reasonable results after only a small number of iterations.

Whenever λ is fixed to a particular positive value $\lambda^* > 0$, the problem (12) simplifies to an unconstrained quadratic program (QP) in ν of the form $\min_{\nu \in \mathbb{R}^{\bar{d}}} \tilde{C}_x(\lambda^*, \nu)$. On the other hand, if ν is fixed to $\nu^* \in \mathbb{R}^{\bar{d}}$, the dual objective takes the form

$$\min_{\lambda \in \mathbb{R}_{>0}} \tilde{C}_x(\lambda, \nu^*) = \min_{\lambda \in \mathbb{R}_{>0}} \frac{c_1}{\lambda} + c_2 \lambda + c_3 \quad (13)$$

with the constants

$$\begin{aligned} c_1 &= \frac{1}{4} \begin{bmatrix} \Lambda^\top \nu^* \\ -1 \end{bmatrix}^\top K_{\mathbb{X}\mathbb{X}} \begin{bmatrix} \Lambda^\top \nu^* \\ -1 \end{bmatrix}, \\ c_2 &= \Gamma^2, \quad c_3 = y^\top \nu^* + \bar{\delta} \|\nu^*\|_1 \end{aligned} \quad (14)$$

and $K_{\mathbb{X}\mathbb{X}} \succeq 0$. We have $\frac{\partial \tilde{C}_x(\lambda, \nu^*)}{\partial \lambda} = \frac{-c_1}{\lambda^2} + c_2$ which gives the candidate solution $\lambda^* = \sqrt{c_1/c_2}$ for (13). We have $c_2 > 0$ and $c_1 \geq 0$ for any $x \in \Omega$, and $c_1 > 0$ for any $x \in \Omega \setminus X$. Furthermore, λ^* is indeed a minimizer of (13) for $x \in \Omega \setminus X$ since its curvature is positive, i.e., $\frac{\partial^2 \tilde{C}_x(\lambda^*, \nu^*)}{\partial \lambda^2} > 0$. In closed-form, λ^* takes the following form.

$$\lambda^* = \frac{\sqrt{\nu^{*\top} \Lambda K_{XX} \Lambda^\top \nu^* - 2(\Lambda K_{XX})^\top \nu^* + k(x, x)}}{2\Gamma} \quad (15)$$

Note that $\lambda^* = 0$ is only possible if $\begin{bmatrix} \Lambda^\top \nu^* \\ -1 \end{bmatrix}$ is in the nullspace of matrix $K_{\mathbb{X}\mathbb{X}}$, which is only possible if $x \in X$. In this case, after fixing $\lambda^* = 0$, the problem to solve for ν reduces to

$$\min_{\nu} y^\top \nu + \bar{\delta} \|\nu\|_1 \quad \text{s.t.} \quad \begin{bmatrix} \Lambda^\top \nu \\ -1 \end{bmatrix} \in \text{Null}(K_{\mathbb{X}\mathbb{X}}).$$

We formulate the alternating optimization algorithm for a maximum number of iterations L and a termination threshold ϵ in the following way.

Algorithm 1: Alternating minimization

Result: Upper bound $\tilde{C}(x)$ of the ground-truth at point x

Input: $x, \lambda_0, L, \epsilon$

$\lambda_0^* = \lambda_0$

$k = 0$

do

$\nu^* = \arg \min_{\nu \in \mathbb{R}^{\bar{d}}} \tilde{C}_x(\lambda_k^*, \nu)$

$\lambda_{k+1}^* = \frac{\sqrt{\nu^{*\top} \Lambda K_{XX} \Lambda^\top \nu^* - 2(\Lambda K_{XX})^\top \nu^* + k(x, x)}}{2\Gamma}$

$k = k + 1$

while $k < L$ and $|\lambda_k^* - \lambda_{k-1}^*| > \epsilon$;

$\tilde{C}(x) = \tilde{C}_x(\lambda_k^*, \nu^*)$

Remark 6 (Numerical properties). Recall the convex dual objective function (12). Since the non-differentiable term $\|\nu\|_1$ is separable and the remainder of the objective is differentiable, a tuple (ν^*, λ^*) that simultaneously minimizes both sub-problems also necessarily minimizes the whole objective (12). For non-asymptotic sublinear convergence rates of alternating minimization algorithms applied to convex programs, the reader is referred to the work [42].

V. ON THE CONNECTIONS WITH GAUSSIAN PROCESSES

Before proceeding to a first numerical example, we contrast our uncertainty quantification technique with Gaussian processes (GPs). By putting them into perspective, we hope to improve the understanding of the theory developed herein.

In the Bayesian framework of GPs, kernels are used to parametrize the covariance between random variables in the input space Ω [43]. Deeper links also exist between the Hilbert space associated with such stochastic processes and the RKHS \mathcal{H} corresponding to their kernels [44], in that there exists an isometric isomorphism connecting both spaces. Moreover, it is known that even though the mean function of a GP does belong to \mathcal{H} , its sample paths almost surely do not if $\dim(\mathcal{H}) = \infty$. Nonetheless, the same paths can belong to another RKHS—with probability one (see [45] for a comprehensive discussion on the topic). The latter phenomenon is known as Driscoll's zero-one law. Although they unveil fundamental properties of the GP framework, deriving practical guidelines from these results requires care as intuition might lead to wrong conclusions when examining infinite-dimensional spaces.

The variance of a Gaussian process is a form of uncertainty quantification against outputs y drawn from the distribution

conditioned on a query input x . Some other GP works are more aligned with our setting and consider a ground-truth map, either assumed to be a GP sample path [46] or to belong to the corresponding kernel RKHS [47], [48]. The latter works bound the difference between the GP mean and the ground-truth values, making use of the GP standard deviation times a uniform scaling term. As for noise models, or marginal likelihood in statistical terms, the most widely adopted forms are the Gaussian or sub-Gaussian formats. The primary motivation behind this choice is analytical tractability since the GP variance does not admit a closed-form expression and has to be numerically approximated in case other noise models are used (e.g. heteroscedastic Gaussian, log-Gaussian, Bernoulli) [43, Chapter 9].

Whether one should opt for the theory developed herein or for Gaussian processes-based techniques is truly a question of model selection. If used to analyze kernel models of dynamical systems, our approach would allow for the use of robust analysis and control tools since worst-case effects and distances can be computed. On the other hand, carrying out modeling through Gaussian processes requires users to use stochastic control theory [49]. As a result, the final yield should also be taken into account when choosing a technique. Do probabilistic inequalities suffice or does my application require deterministic certification? In our view, the deterministic and the stochastic frameworks have their own merits and the user should judge which of them is more adequate to tackle the problem at hand. Despite the differences between the standing assumptions of the methods, we provide the reader with a comparison between the bounds developed herein and a popular GP alternative in Section VI.

VI. NUMERICAL EXAMPLES

The methods developed in the previous sections are now employed in three distinct scenarios: a function bounding task, an optimization problem with an unknown constraint, and a control certification procedure.¹

Example 1: Consider the function below, which represents the first component update map of a Hénon chaotic attractor with an additional sinusoidal forcing term

$$f^*(z_1, z_2) = 1 - az_1^2 + z_2 + b \sin(cz_2) \quad (16)$$

The parameters are $a = 0.8$, $b = 8$ and $c = 0.8$, and its domain is the box $\Omega = [-10 \ 10] \times [-10 \ 10]$. A squared-exponential kernel $k(x, x') = \exp\left(-\frac{\|x-x'\|^2}{2\ell^2}\right)$ with $x = [z_1 \ z_2]$ was chosen for our experiments with lengthscale $\ell = 5$, which was empirically estimated by gridding the search-space and performing posterior validation. Γ was obtained through the procedure described in Appendix A with a final value of $\Gamma = 1200$. $d = 100$ samples were collected using two strategies: inputs lying in an equidistant grid, and inputs being drawn randomly from a uniform distribution. Noise was sampled uniformly throughout the tests with $\bar{\delta} = 1$ and $\bar{\delta} = 5$.

The obtained optimal upper bound $C(x)$ is displayed in Figure 2 along with the ground-truth function f^* . Consider the

scenarios where $\bar{\delta} = 1$. Whereas the $C(x)$ surface is overall tight for the grid-based dataset, with an average distance of 3.01 to the latent function, randomized data yielded a less regular bound with an average distance of 8.02. These numbers were increased respectively to 9.57 and 18.97 when the noise levels were risen to $\bar{\delta} = 5$. The plots illustrate the disadvantages of relying on completely randomized input locations, which degraded especially the borders of $C(x)$. An equidistant grid of points is highly favorable since it not only fills the domain well, but also ensures a minimum separation distance so that no two inputs are too close to cause numerical problems when handling the kernel matrix K_{XX} .

The $f^*(z_1, z_2)$ map was then sliced at $z_1 = -3$ and the entire envelope $B(x) \leq x \leq C(x)$ was computed. This was compared to the sub-optimal bounds given in Proposition 3 for a kernel ridge regression (KRR) model. The two previous datasets with $\bar{\delta} = 1$ were used and the obtained results are displayed in Figure 3. As can be seen from the plots, the optimal approach yielded tighter uncertainty intervals that were always within the sub-optimal ones. Moreover, whereas the average width of the blue envelope was 8.93 and 18.96 respectively in the grid and random cases, the green envelope displayed average widths of 21.13 and 34.62.

Next, we consider the Gaussian process bounds proposed in [47, Lemma 3] (see also the closely related works [50], [51]) and analyze how they compare to the proposed robust ones. Overloading notation for the sake of clarity, these bounds have the form

$$|\mu(x) - f^*(x)| \leq \beta \sigma(x) \quad (17a)$$

$$\text{with } \beta = \Gamma + 4\lambda\sqrt{\gamma + 1 + \ln(1/\bar{\delta})}, \quad (17b)$$

where $\mu(x)$ is the GP mean, $\sigma(x)$ is its standard deviation, λ is the strength of the sub-Gaussian noise, γ is the maximum information capacity for a fixed number of samples, and $1 - \delta$ is the confidence of the inequality. For a detailed explanation of how the experiment was set up, the reader is referred to Appendix E. The data, $d = 100$ samples, corrupted by the same noise realizations were used throughout the tests for all methods. Two parameters were then varied to understand how sensitive each method is to them: the RKHS norm estimate Γ and the noise bound $\bar{\delta}$, which were increased by a factor of 1, 1.5, and 2. Detailed results can be found in Appendix E, Tables I and II. The outcomes in all 18 different scenarios were unanimous in ranking the optimal bounds as the tightest method, followed by the sub-optimal ones, and finally the GP approach. Indeed, the GP bounds always yielded average widths at least one order of magnitude greater than the optimal deterministic ones. We attribute this difference especially to the direct product between of Γ and $\sigma(x)$ in (17), which causes them to be particularly sensitive to norm over-approximations. This effect is dampened in (9) due to the interaction with $\tilde{\Delta}$ (see the derivation in Appendix C-B).

Example 2: The next numerical experiment involves the ground-truth (16) as an unknown constraint for a static problem (data-driven optimization with unknown constraints is typical in the field of real-time optimization [33]). Consider

¹The code to reproduce our results is available at <https://github.com/PREDICT-EPFL/opt-rkhs-bounds>.

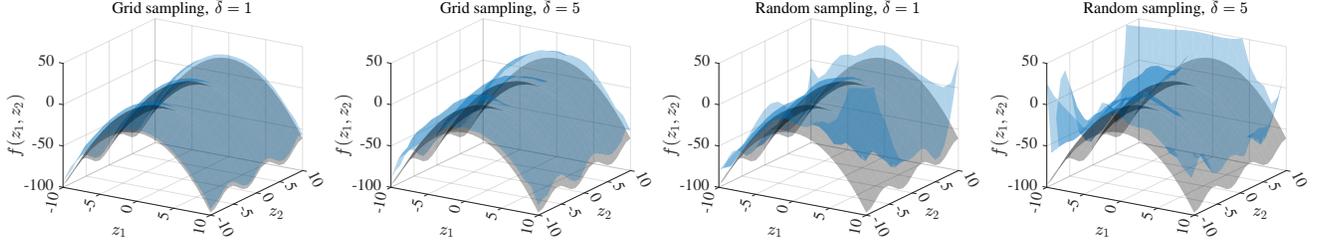


Fig. 2. The ground-truth (black) and the upper optimal bound $C(x)$ (blue) with 100 data-points. Two noise levels are considered, $\bar{\delta} = 1$ and $\bar{\delta} = 5$, and two sampling strategies, an equidistant grid and random uniform sampling.

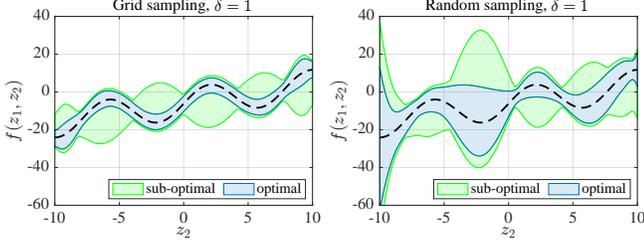


Fig. 3. A comparison between the optimal bounds (blue) and the closed-form sub-optimal ones (green). The 2D ground-truth was sliced at $z_1 = -3$ and is shown in dashed black.

the following formulation

$$\min_{z \in \mathbb{R}^2} (z_1 - 1)^2 + (z_2 - 5)^2 \quad (18a)$$

$$\text{subj. to } f^*(z) \leq -10 \quad (18b)$$

where the function $f^*(z)$ that maps the decision variables to the constraint is not explicitly known, but can be measured. Samples were used to establish an upper bound $C(z)$ for $f^*(z)$, hence providing an inner-approximation for the real feasible set. We considered the cases of having 64, 81 and 100 evaluations of $f^*(z)$ affected by noise with $\bar{\delta} = 1$ and, once more, the data were collected by means of a uniform random distribution and an equidistant grid. In the approximate optimization problems, the original constraint (18b) was replaced by $C(z) \leq -10$. Optimizers z^* were computed by gridding the domain, and the results along with the estimated feasible sets (shaded areas) are shown in Figure 4. Notice how in some instances the set of feasible decisions is not connected. Thanks to Proposition 2, the addition of new data-points can only relax the approximate formulation, hence reducing the found minimum. Indeed, the obtained solutions for the approximate problems were 13.21, 11.36 and 10.96, respectively with 64, 81 and 100 samples taken randomly. When employing a grid, the figures were 10.67, 8.48 and 7.67. The solution of the real problem, i.e., the one with the ground-truth constraint, is 5.69.

Example 3: Finally, we verify if a sequence of control actions obtained by means of a certainty equivalence approach will or will not lead to the real system violating constraints. In this scenario, previous examples in the form of control and state trajectories are exploited to build the necessary datasets.

Let us consider a continuous stirred-tank reactor (CSTR)

described by the differential equations

$$\dot{c}_A(t) = u(t)(c_{A0} - c_A(t)) - \rho_1 c_A(t) - \rho_3 c_A^2(t) \quad (19a)$$

$$\dot{c}_B(t) = -u(t)c_B(t) + \rho_1 c_A(t) - \rho_2 c_B^2(t) \quad (19b)$$

c_A and c_B denote respectively the concentrations of cyclopentadiene and cyclopentenol, whereas u represents the feed inflow of cyclopentadiene. The parameters are $\rho_1 = \rho_2 = 4.1 \times 10^{-3} \text{ h}^{-1}$, $\rho_3 = 6.3 \times 10^{-4} \text{ h}^{-1}$, $c_{A0} = 5.1 \text{ mol/l}$. The system is subject to the constraints $1 \leq c_A \leq 3$, $0.5 \leq c_B \leq 2$, $3 \leq u \leq 25$, and is sampled at a rate of $1/30 \text{ Hz}$. In order to steer the CSTR states toward $c_A^{\text{ref}} = 2.14$, $c_B^{\text{ref}} = 1.09$, an optimal control problem (OCP) based on KRR models was formulated and solved. The approach featured no uncertainty quantification, i.e., it relied solely on certainty equivalence.

To verify if OCP control actions would not lead to the true system violating constraints, the tools developed in Section III were employed. The certification problem was broken down into several steps: the 1-step ahead analysis, the 2-step ahead analysis, and so on. The associated datasets $\{(x_i, y_i)\}_{i=1}^d$ were composed of initial states and sequences of control actions to form x_i , and the final state to form y_i (this multi-step approach is the same as the one explained in [30, Sec. 4]). The squared-exponential kernel was used throughout the whole process and the various lengthscales were selected through a 5-fold cross-validation process based on 400 samples. The same batch of data were exploited to estimate the different RKHS norms Γ following the procedure of Appendix A. The obtained lower estimates were then augmented to account for possible unseen complexity. A different dataset was gathered to compute the bounds by starting the system at 800 initial conditions and solving OCP from those locations. We highlight that, as there are two states and one control variable, the domain of the ground-truth mapping the initial condition to the 8-th step ahead state has dimension 10, hence justifying the need for a large dataset. The noise affecting the measurements was drawn uniformly from the interval -0.05 to 0.05 , and a bound of $\bar{\delta} = 0.06$ was used.

The two types of bounds were then computed defining an interval per state and, thus, a ‘‘bounding box’’ for each step. These are then guaranteed to contain the true system evolution, the ground-truth, as illustrated in Figure 5. After visually inspecting the phase portraits, one sees how conservative the sub-optimal method was: the average area of the sub-optimal boxes was 0.1780, and 0.0741 in the optimal case. In addition to it, one bounding box around the trajectory that starts at the bottom right corner of the plots extends below the $c_B \geq 0.5$

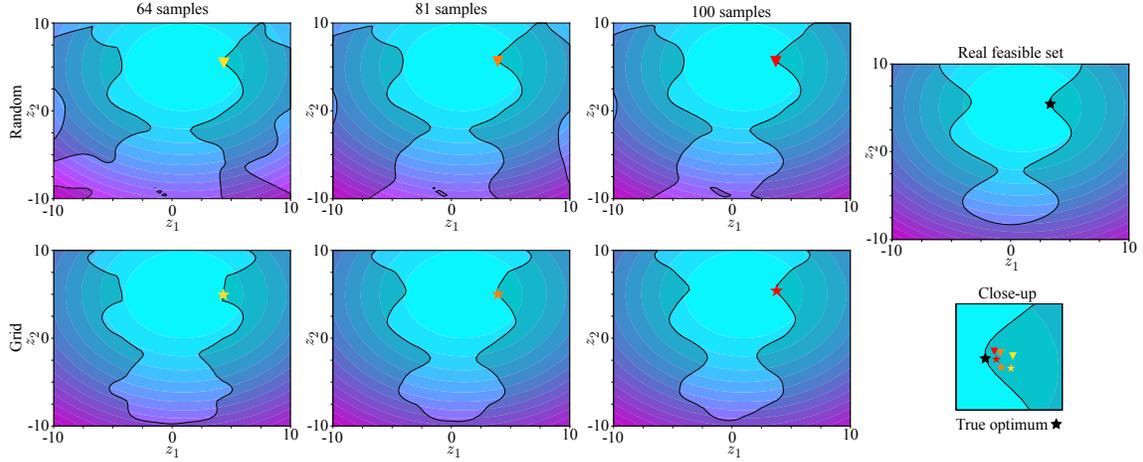


Fig. 4. Solutions and feasible sets (shaded areas) for problem (18) with 64, 81 and 100 samples of $f^*(z)$. Top row: samples drawn uniformly. Bottom row: samples on an equidistant grid. The true feasible set and optimal solution are shown on the right.

constraint. A time-domain view of the situation is shown in Figure 6, where at time-step 1 the lower-bound violates the aforementioned constraint in the top plot, but not in the bottom one. As a result, the OCP control sequence could not be certified by the sub-optimal approach, but could by means of the optimal one.

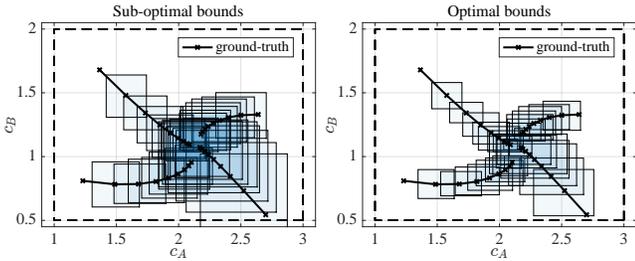


Fig. 5. Phase portraits of the CSTR system under the same control inputs, but with different uncertainty quantification techniques. Constraints are represented by the dashed lines.

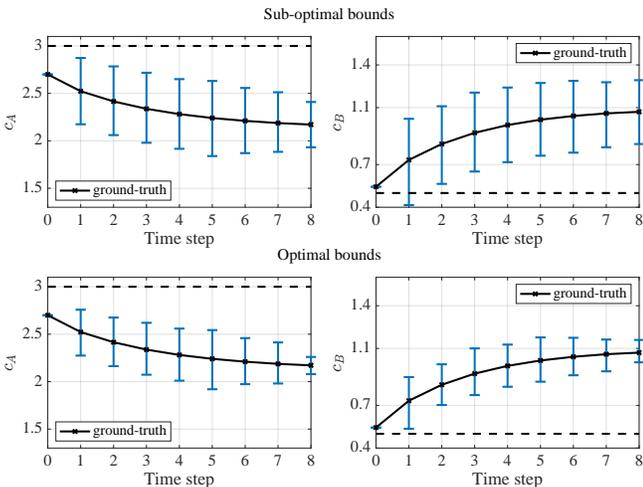


Fig. 6. State trajectories of the CSTR system under the same control inputs, but with different uncertainty quantification techniques. Constraints are represented by the dashed lines.

VII. FINAL REMARKS AND FUTURE DIRECTIONS

We investigated the uncertainty quantification problem associated with evaluations of an unknown function that belongs to a possibly infinite-dimensional reproducing-kernel Hilbert space. Optimal robust bounds were derived by exploiting a finite set of samples and an estimate of the ground-truth function complexity as measured by its norm. Several formulations were then analyzed: a primal finite-dimensional program, one possible dual form, as well as closed-form sub-optimal solutions centered around pre-specified kernel models. When considering the optimal alternatives, it was shown how the addition of new data can only shrink the uncertainty envelope everywhere.

Future research could focus on the following topics. Firstly, the developed theory could be generalized to accept uncertain inputs, thus allowing for uncertainty propagation in multi-stage problems. Additionally, resampling techniques could be used to construct sparse representations or to confer a desired geometrical property on the input points, enabling fast evaluation of the bounds. Exploring further estimation techniques for Γ and $\bar{\delta}$, especially joint estimation, could be of interest for practical application of the approach. Finally, the developed finite-sample bounds could give support to the area of data-driven optimization under unknown constraints or objectives by establishing formal feasibility or performance guarantees.

APPENDIX A ESTIMATING RKHS COMPLEXITY FROM DATA

We consider an unknown map $f \in \mathcal{H}$ and a set of samples $D = \{(x_i, f(x_i))\}_{i=1}^d$. Using the shorthand $f_X = [f(x_1) \dots f(x_d)]^\top$, we have that

$$\hat{\Gamma} := \sqrt{f_X^\top K_{XX}^{-1} f_X} \leq \|f\|_{\mathcal{H}} \quad (20)$$

for any number of samples $d \in \mathbb{N}$ due to the optimal recovery property [4]. Moreover, the decomposition used in the proof of Proposition 2 shows that the quantity $\hat{\Gamma}$ can only

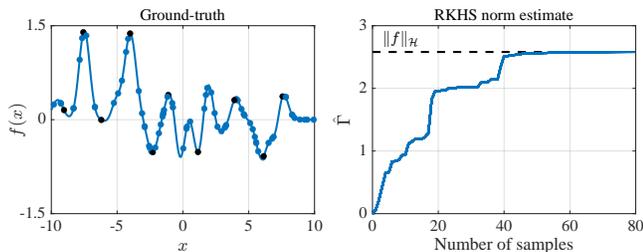


Fig. 7. Estimating the RKHS norm using randomly sampled data (all circles). The quadratic form $\hat{\Gamma}$ for the random samples is shown on the right plot. If one sampled only the black subset of the data, the corresponding $\hat{\Gamma}$ would capture over 90% of the total complexity, i.e., $\hat{\Gamma}/\|f\|_{\mathcal{H}} > 0.9$.

increase with the addition of new data. Since $\|f\|_{\mathcal{H}}$ is the least upper bound for it, then this quantity can be used as an efficient lower estimate for the RKHS norm. An example is shown in Figure 7 for an f composed of 25 squared-exponential kernel functions from which 80 samples were drawn uniformly (left). The corresponding values for $\hat{\Gamma}$ for an increasing number of data are also reported (right). After around 40 samples, essentially all of the RKHS complexity had already been captured. Moreover, by sampling only the peaks and valleys indicated by the black subset of the data-points, one could retrieve over 90% of the final norm. In a practical situation, expert knowledge should be elicited to augment $\hat{\Gamma}$ through a safety factor and hopefully transform it into an upper bound $\Gamma \geq \|f\|_{\mathcal{H}}$. Note however that no hard guarantees are offered—a situation similar to estimating Lipschitz constants purely from scattered observations. Finally, in case the outputs are corrupted by measurement noise, it is possible to quantify its worst-case effect on the estimation process [13].

APPENDIX B THE DATA-SELECTION MATRIX

Recall that n_1, n_2, \dots, n_d are the number of outputs available at the input locations x_1, x_2, \dots, x_d . Λ has size $(\sum_i n_i) \times d$ and is defined as

$$\Lambda := \begin{bmatrix} \mathbf{1}_{n_1} & \mathbf{0}_{n_1} & \mathbf{0}_{n_1} & \cdots & \mathbf{0}_{n_1} \\ \mathbf{0}_{n_2} & \mathbf{1}_{n_2} & \mathbf{0}_{n_2} & \cdots & \mathbf{0}_{n_2} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \mathbf{0}_{n_d} & \mathbf{0}_{n_d} & \mathbf{0}_{n_d} & \cdots & \mathbf{1}_{n_d} \end{bmatrix} \quad (21)$$

where $\mathbf{1}_{n_i}$ and $\mathbf{0}_{n_i}$ are respectively column vectors of ones and zeros of size n_i . If only a single output is available at every input, Λ simplifies to an identity matrix.

APPENDIX C DERIVATIONS

A. Proof of Theorem 1

Let $\mathbb{X} := X \cup \{x\}$ and define the finite-dimensional subspace $\mathcal{H}^{\parallel} = \{f \in \mathcal{H} : f \in \text{span}(k(x_i, \cdot), x_i \in \mathbb{X})\}$. Furthermore, let $\mathcal{H}^{\perp} = \{g \in \mathcal{H} : \langle g, f^{\parallel} \rangle_{\mathcal{H}} = 0, \forall f^{\parallel} \in \mathcal{H}^{\parallel}\}$ be the orthogonal complement of \mathcal{H}^{\parallel} . Then, we have $\mathcal{H} = \mathcal{H}^{\parallel} \oplus \mathcal{H}^{\perp}$ and for all $f \in \mathcal{H}$, $\exists f^{\parallel} \in \mathcal{H}^{\parallel}, f^{\perp} \in \mathcal{H}^{\perp} : f = f^{\parallel} + f^{\perp}$. By

employing the latter decomposition and using the reproducing property, we can reformulate $\mathbb{P}0$ in terms of \mathcal{H}^{\parallel} and \mathcal{H}^{\perp} as

$$\sup_{\substack{f^{\parallel} \in \mathcal{H}^{\parallel} \\ f^{\perp} \in \mathcal{H}^{\perp}}} \left\{ \langle f^{\parallel} + f^{\perp}, k(x, \cdot) \rangle_{\mathcal{H}} : \left\| f^{\parallel} + f^{\perp} \right\|_{\mathcal{H}}^2 \leq \Gamma^2, \left\| (f^{\parallel} + f^{\perp})_X - y \right\|_{\infty} \leq \bar{\delta} \right\} \quad (22)$$

$$\stackrel{(i)}{=} \sup_{\substack{f^{\parallel} \in \mathcal{H}^{\parallel} \\ f^{\perp} \in \mathcal{H}^{\perp}}} \left\{ f^{\parallel}(x) : \left\| f^{\parallel} \right\|_{\mathcal{H}}^2 + \left\| f^{\perp} \right\|_{\mathcal{H}}^2 \leq \Gamma^2, \left\| f^{\parallel}_X - y \right\|_{\infty} \leq \bar{\delta} \right\} \quad (23)$$

$$\stackrel{(ii)}{=} \sup_{f^{\parallel} \in \mathcal{H}^{\parallel}} \left\{ f^{\parallel}(x) : \left\| f^{\parallel} \right\|_{\mathcal{H}}^2 \leq \Gamma^2, \left\| f^{\parallel}_X - y \right\|_{\infty} \leq \bar{\delta} \right\} \quad (24)$$

In (i), the f^{\perp} component vanished from the cost and from the last constraint due to orthogonality w.r.t. $k(x_i, \cdot) \in \mathcal{H}^{\parallel}$ for any $x_i \in \mathbb{X}$; moreover, the Pythagorean relation $\|f\|_{\mathcal{H}}^2 = \|f^{\parallel}\|_{\mathcal{H}}^2 + \|f^{\perp}\|_{\mathcal{H}}^2$ was also used. To arrive at the second equality (ii), one only has to note that the objective is insensitive to f^{\perp} and that any $f^{\perp} \neq 0_{\mathcal{H}}$ would tighten the first constraint.

The attainment of the supremum is addressed next. Consider (24) and denote the members of \mathcal{H}^{\parallel} simply as f . $\|f\|_{\mathcal{H}}^2 \leq \Gamma^2$ is a closed and bounded constraint as it is the sublevel set of a norm. We transform $\|f_X - y\|_{\infty} \leq \bar{\delta}$ into $|f(x_i) - y_{i,j}| \leq \bar{\delta}$, $i = 1, \dots, d, j = 1, \dots, n_i$. Sets of the form $\{a \in \mathbb{R} : |a| \leq b\}$ are clearly closed in \mathbb{R} , hence $\{f(x_i) \in \mathbb{R} : |f(x_i) - y_{i,j}| \leq \bar{\delta}, \forall i, j\}$ is also closed. For any x_i , the evaluation functional $L_{x_i}(f) = f(x_i)$ is a linear operator and thus pre-images of closed sets are also closed. Consequently, $\{f \in \mathcal{H}^{\parallel} : |f(x_i) - y_{i,j}| \leq \bar{\delta}, \forall i, j\}$ is closed in \mathcal{H}^{\parallel} . The intersection of a finite number of closed sets is necessarily closed, thus all constraint present in (24) define a closed feasible set. Since \mathcal{H}^{\parallel} is finite-dimensional, any closed and bounded subset of it is compact (Heine–Borel); therefore, the continuous objective $L_x(f) = f(x)$ in (24) attains a maximum by the Weierstrass extreme value theorem.

Finally, we establish the connection between $\mathbb{P}0$ and $\mathbb{P}1$. From the above arguments, an optimizer for $\mathbb{P}0$ must lie in \mathcal{H}^{\parallel} . The members $f \in \mathcal{H}^{\parallel}$ have the form $f(z) = \alpha^{\top} K_{\mathbb{X}\mathbb{X}} z$, being defined by the α weights. Due to the positive-definiteness of k , there exists a bijective map between outputs at the \mathbb{X} locations $f_{\mathbb{X}} = [f(x_1) \dots f(x_d) f(x)]^{\top}$ and the weights α , namely $\alpha = K_{\mathbb{X}\mathbb{X}}^{-1} f_{\mathbb{X}}$. $K_{\mathbb{X}\mathbb{X}}$ denotes the kernel matrix associated with the set $\mathbb{X} = X \cup \{x\}$. Consequently, optimizing over $f \in \mathcal{H}^{\parallel}$ is equivalent to optimizing over $[f(x_1) \dots f(x_d) f(x)]^{\top} =: [c^{\top} \ c_x]^{\top}$. The bounded norm condition can be recast as $\|f\|_{\mathcal{H}}^2 = \langle f, f \rangle_{\mathcal{H}} = \alpha^{\top} K_{\mathbb{X}\mathbb{X}} \alpha = [c^{\top} \ c_x] K_{\mathbb{X}\mathbb{X}}^{-1} [c^{\top} \ c_x]^{\top}$. The remaining constraint and the objective are straightforward. Noting that this reformulation is valid for any $x \in \Omega$ concludes the proof. \square

B. Proof of Proposition 3

For any given $s(x) = \alpha^\top K_{Xx}$, we have

$$\begin{aligned} & |f^*(x) - s(x)| \\ &= |f^*(x) - \tilde{s}(x) + \tilde{s}(x) - s(x)| \end{aligned} \quad (25)$$

$$\leq |f^*(x) - (f_X^* + \delta_X)K_{XX}^{-1}K_{Xx}| + |\tilde{s}(x) - s(x)| \quad (26)$$

$$\leq |f^*(x) - \bar{s}(x)| + \bar{\delta} \|K_{XX}^{-1}K_{Xx}\|_1 + |\tilde{s}(x) - s(x)| \quad (27)$$

$$\leq P(x) \sqrt{\Gamma^2 - \|\bar{s}\|_{\mathcal{H}}^2} + \bar{\delta} \|K_{XX}^{-1}K_{Xx}\|_1 + |\tilde{s}(x) - s(x)| \quad (28)$$

$$\leq P(x) \sqrt{\Gamma^2 + \Delta - \|\bar{s}\|_{\mathcal{H}}^2} + \bar{\delta} \|K_{XX}^{-1}K_{Xx}\|_1 + |\tilde{s}(x) - s(x)| \quad (29)$$

with f_X^* being the vector of true function values at the sample locations in X and δ_X the vector of additive measurement noise for the samples y . (26) follows from the triangle inequality and the additive noise property of y . Using the triangle inequality again, we arrive at (27), where \bar{s} denotes the noise-free interpolant of f_X^* . The noise-free interpolation error bound gives the estimation in the first term of (28), while (29) follows from [13, Lemma 1], with $\Delta = \max_{\|\delta\|_\infty \leq \bar{\delta}} (-\delta^\top K_{XX}^{-1} \delta + 2y^\top K_{XX}^{-1} \delta)$. A standard dualization procedure as the one presented in Appendix C-C leads to the dual problem

$$\min_{\nu \in \mathbb{R}^d} \frac{1}{4} \nu^\top K_{XX} \nu + \nu^\top y + \bar{\delta} \|\nu\|_1 + y^\top K_{XX}^{-1} y \quad (30)$$

for Δ . Notice that the last term in (30) is constant and the same as the squared interpolant norm $\|\bar{s}\|_{\mathcal{H}}^2$. Therefore, these terms cancel in (29) and we are left with

$$\begin{aligned} |f^*(x) - s(x)| &\leq P(x) \sqrt{\Gamma^2 + \tilde{\Delta}} + \bar{\delta} \|K_{XX}^{-1}K_{Xx}\|_1 \\ &\quad + |\tilde{s}(x) - s(x)| \end{aligned} \quad (31)$$

where $\tilde{\Delta}$ represents (30) without the constant term.

C. The Lagrange dual problem

Consider the case $x \notin X$. Let $z := [c^\top \quad c_x]^\top$, $a := [\mathbf{0}^\top \quad 1]^\top$, $A := [\mathbf{I} \quad \mathbf{0}]$. The Lagrangian of $\mathbb{P}1$ is

$$\begin{aligned} \mathcal{L}(z, \lambda, \beta, \gamma) &= a^\top z - \lambda(z^\top K_{XX}^{-1} z - \Gamma^2) \\ &\quad - \beta^\top (\Lambda A z - y - \bar{\delta} \mathbf{1}) - \gamma^\top (y - \Lambda A z - \bar{\delta} \mathbf{1}) \end{aligned} \quad (32)$$

where K_{XX} denotes the kernel matrix evaluated at $X \cup \{x\}$. Suppose $\lambda > 0$. Computing $\nabla_z \mathcal{L}(z^*) = 0$ leads to

$$z^* = -\frac{1}{2\lambda} K_{XX} (A^\top \Lambda^\top (\beta - \gamma) - a).$$

Defining the auxiliary variable $\nu = \beta - \gamma$, and substituting z^* into (32) gives the dual objective

$$\begin{aligned} g(\lambda, \nu) &= \frac{1}{4\lambda} \nu^\top \Lambda A K_{XX} A^\top \Lambda^\top \nu + \left(y - \frac{1}{2\lambda} \Lambda A K_{XX} a \right)^\top \nu \\ &\quad + \bar{\delta} \|\nu\|_1 + \frac{1}{4\lambda} a^\top K_{XX} a + \lambda \Gamma^2 \end{aligned} \quad (33)$$

$$\begin{aligned} &= \frac{1}{4\lambda} \nu^\top \Lambda K_{XX} \Lambda^\top \nu + \left(y - \frac{1}{2\lambda} \Lambda K_{Xx} \right)^\top \nu \\ &\quad + \bar{\delta} \|\nu\|_1 + \frac{1}{4\lambda} k(x, x) + \lambda \Gamma^2 \end{aligned} \quad (34)$$

where in the second equality the matrix K_{XX} was expanded and the resulting terms were reorganized. Since $\beta, \gamma \in \mathbb{R}_{\geq 0}^d$ and $\nu = \beta - \gamma$ then ν is unconstrained.

Now if $\lambda = 0$, the Lagrangian (32) simplifies to $\mathcal{L}(z, \nu) = (a - A^\top \Lambda^\top \nu)^\top z + \nu^\top y + \bar{\delta} \|\nu\|_1$, which is linear in z . Its supremum w.r.t. z is only finite if $a = A^\top \Lambda^\top \nu$. Recalling the definitions of a , A and Λ , one can see that $\nexists \nu$ that could satisfy the latter condition. Therefore, $\lambda = 0 \implies \sup_z \mathcal{L}(z, \lambda, \nu) = +\infty$, meaning that the dual problem is infeasible. As a conclusion, the Lagrangian dual of $\mathbb{P}1$ in (6) is precisely $\mathbb{D}1$ in (12).

Next, consider the case $x \in X$, $x = x_i$. The objective of $\mathbb{P}1'$ can be written as $a^\top c$ with $a_i = 1$ and $a_n = 0, n \neq i$. When deriving its Lagrangian, one obtains again (32) with the simplifications: $z \leftarrow c$, $K_{XX} \leftarrow K_{XX}$ and $A \leftarrow \mathbf{I}$. We proceed by analyzing the two scenarios for λ as before. If $\lambda > 0$, the previous derivations apply, leading to the same the quadratic-over-linear objective (34). However, if $\lambda = 0$, the Lagrangian becomes $\mathcal{L}(z, \nu) = (a - \Lambda^\top \nu)^\top z + \nu^\top y + \bar{\delta} \|\nu\|_1$, whose supremum w.r.t. z is only finite if $a = \Lambda^\top \nu$. In contrast with the previous paragraph, this condition now can be satisfied. It is equivalent to $\nu_{i,1} + \dots + \nu_{i,n_i} = 1$, where the variables are all the multipliers associated with the i -th input location x_i . The resulting expression can be minimized analytically, yielding the minimum $\min_j y_{i,j} + \bar{\delta}$, i.e., the smallest output available at x_i augmented by the noise bound. Finally, we conclude that the dual objective for $\mathbb{P}1'$ is

$$g(\lambda, \nu) = \begin{cases} (34), & \text{if } \lambda > 0 \\ \min_j y_{i,j} + \bar{\delta}, & \text{if } \lambda = 0 \end{cases} \quad (35)$$

As a last observation, a dual problem can also be derived for (7), calculating the lower part of the envelope. The formulation is analogous to (12), assuming the form

$$\begin{aligned} \max_{\nu \in \mathbb{R}^d, \lambda > 0} & -\frac{1}{4\lambda} \nu^\top \Lambda K_{XX} \Lambda^\top \nu - \left(y + \frac{1}{2\lambda} \Lambda K_{Xx} \right)^\top \nu \\ & - \bar{\delta} \|\nu\|_1 - \frac{1}{4\lambda} k(x, x) - \lambda \Gamma^2 \end{aligned} \quad (36)$$

Note that these are distinct objectives, not merely opposites. Therefore, two problems have to be solved to fully quantify the ground-truth uncertainty.

APPENDIX D

A BLOCK MATRIX IDENTITY

Let $A \in \mathbb{R}^{d \times d}$ be invertible, $B \in \mathbb{R}^d$ and $c \in \mathbb{R}$. The following identity holds

$$\begin{bmatrix} A & B \\ B^\top & c \end{bmatrix}^{-1} = \begin{bmatrix} A^{-1} + \frac{1}{d} A^{-1} B B^\top A^{-1} & -\frac{1}{d} A^{-1} B \\ -\frac{1}{d} B^\top A^{-1} & \frac{1}{d} \end{bmatrix} \quad (37)$$

where $d = c - B^\top A^{-1} B$.

APPENDIX E

GP COMPARISON: SETTINGS AND RESULTS

In order to compare the GP uncertainty bounds (17) to their deterministic counterparts, the following approach was

adopted. First, a lower bound for the maximum information gain γ was used since the problem of exactly computing such a quantity is in general NP-hard [52]. Note how this decision favors the GP bounds by shrinking them. The chosen lower bound was the information gain of our inputs X , which in our zero-mean Gaussian noise setting with variance λ^2 is $\frac{1}{2} \ln(\det(I + \lambda^{-2} K_{XX}))$ [52]. As for the noise realizations, we proceeded as follows. Starting from our hard noise limit $\bar{\delta}$, we considered a zero-mean Gaussian distribution with variance such that its samples would lie in the $[-\bar{\delta}, \bar{\delta}]$ band with confidence 0.99, i.e., a standard deviation of $\lambda = \frac{\bar{\delta}}{2.58}$. The noise was then drawn from the normal distribution and clipped to the interval $[-\bar{\delta}, \bar{\delta}]$ to fulfill Assumption 1. Finally, the probabilistic inequality (17) was evaluated for a final confidence of 99%. The obtained numerical results are shown in Tables I and II.

REFERENCES

- [1] B. Schölkopf, "Causality for machine learning," in *Probabilistic and Causal Inference: The Works of Judea Pearl*, 2022, pp. 765–804.
- [2] F. Cucker and D. X. Zhou, *Learning theory: An approximation theory viewpoint*. Cambridge University Press, 2007, vol. 24.
- [3] R. Schaback, "Error estimates and condition numbers for radial basis function interpolation," *Advances in Computational Mathematics*, vol. 3, no. 3, pp. 251–264, 1995.
- [4] H. Wendland, *Scattered data approximation*. Cambridge university press, 2004, vol. 17.
- [5] A. Iske, *Approximation Theory and Algorithms for Data Analysis*. Springer, 2018.
- [6] M. Belkin, "Approximation beats concentration? an approximation view on inference with smooth radial kernels," in *Conference On Learning Theory*. PMLR, 2018, pp. 1348–1361.
- [7] S. Mei, T. Misiakiewicz, and A. Montanari, "Mean-field theory of two-layers neural networks: dimension-free bounds and kernel limit," in *Conference on Learning Theory*. PMLR, 2019, pp. 2388–2464.
- [8] P. Domingos, "Every model learned by gradient descent is approximately a kernel machine," *arXiv preprint arXiv:2012.00152*, 2020.
- [9] A. De Mathews, J. Hron, M. Rowland, R. Turner, and Z. Ghahramani, "Gaussian process behaviour in wide deep neural networks," in *6th International Conference on Learning Representations, ICLR 2018-Conference Track Proceedings*, 2018.
- [10] B. Schölkopf, A. J. Smola, F. Bach *et al.*, *Learning with kernels: support vector machines, regularization, optimization, and beyond*. MIT press, 2002.
- [11] C. A. Micchelli, Y. Xu, and H. Zhang, "Universal kernels," *Journal of Machine Learning Research*, vol. 7, no. 12, 2006.
- [12] R. Schaback and H. Wendland, "Kernel techniques: from machine learning to meshless methods," *Acta numerica*, vol. 15, p. 543, 2006.
- [13] E. T. Maddalena, P. Scharnhorst, and C. N. Jones, "Deterministic error bounds for kernel-based learning techniques under bounded noise," *Automatica*, vol. 134, p. 109896, 2021.
- [14] T. O'Hagan, "Dicing with the unknown," *Significance*, vol. 1, no. 3, pp. 132–133, 2004.
- [15] A. Chakrabarty, V. Dinh, M. J. Corless, A. E. Rundell, S. H. Żak, and G. T. Buzzard, "Support vector machine informed explicit nonlinear model predictive control using low-discrepancy sequences," *IEEE Transactions on Automatic Control*, vol. 62, no. 1, pp. 135–148, 2016.
- [16] U. Rosolia and F. Borrelli, "Learning model predictive control for iterative tasks. a data-driven control framework," *IEEE Transactions on Automatic Control*, vol. 63, no. 7, pp. 1883–1896, 2017.
- [17] J. Umlauf and S. Hirche, "Feedback linearization based on gaussian processes with event-triggered online learning," *IEEE Transactions on Automatic Control*, vol. 65, no. 10, pp. 4154–4169, 2019.
- [18] T. Chen, M. S. Andersen, L. Ljung, A. Chiuso, and G. Pillonetto, "System identification via sparse multiple kernel-based regularization using sequential convex optimization techniques," *IEEE Transactions on Automatic Control*, vol. 59, no. 11, pp. 2933–2945, 2014.
- [19] 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*, vol. 50, no. 3, pp. 657–682, 2014.
- [20] G. Pin, A. Assalone, M. Lovera, and T. Parisini, "Non-asymptotic kernel-based parametric estimation of continuous-time linear systems," *IEEE Transactions on Automatic Control*, vol. 61, no. 2, pp. 360–373, 2015.
- [21] Y. Fujimoto, I. Maruta, and T. Sugie, "Input design for kernel-based system identification from the viewpoint of frequency response," *IEEE Transactions on Automatic Control*, vol. 63, no. 9, pp. 3075–3082, 2018.
- [22] S. Z. Rizvi, J. M. Velni, F. Abbasi, R. Tóth, and N. Meskin, "State-space LPV model identification using kernelized machine learning," *Automatica*, vol. 88, pp. 38–47, 2018.
- [23] V. Laurain, R. Tóth, D. Piga, and M. A. H. Darwish, "Sparse rkhs estimation via globally convex optimization and its application in LPV-IO identification," *Automatica*, vol. 115, p. 108914, 2020.
- [24] R. S. Risuleo, G. Bottegal, and H. Hjalmarsson, "A nonparametric kernel-based approach to Hammerstein system identification," *Automatica*, vol. 85, pp. 234–247, 2017.
- [25] R. S. Risuleo, F. Lindsten, and H. Hjalmarsson, "Bayesian nonparametric identification of Wiener systems," *Automatica*, vol. 108, p. 108480, 2019.
- [26] M. Lauricella and L. Fagiano, "Set membership identification of linear systems with guaranteed simulation accuracy," *IEEE Transactions on Automatic Control*, vol. 65, no. 12, pp. 5189–5204, 2020.
- [27] F. Berkenkamp, A. P. Schoellig, and A. Krause, "Safe controller optimization for quadrotors with gaussian processes," in *2016 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2016, pp. 491–496.
- [28] A. Lederer, A. Capone, and S. Hirche, "Parameter optimization for learning-based control of control-affine systems," in *Learning for Dynamics and Control*. PMLR, 2020, pp. 465–475.
- [29] J. M. Manzano, D. Limon, D. M. de la Peña, and J.-P. Calliess, "Robust learning-based MPC for nonlinear constrained systems," *Automatica*, vol. 117, p. 108948, 2020.
- [30] E. T. Maddalena, P. Scharnhorst, Y. Jiang, and C. N. Jones, "KPC: Learning-based model predictive control with deterministic guarantees," in *Learning for Dynamics and Control*. PMLR, 2021, pp. 1015–1026.
- [31] C. Knuth, G. Chou, N. Ozay, and D. Berenson, "Planning with learned dynamics: Probabilistic guarantees on safety and reachability via Lipschitz constants," *IEEE Robotics and Automation Letters*, 2021.
- [32] K. P. Wabersich, L. Hewing, A. Carron, and M. N. Zeilinger, "Probabilistic model predictive safety certification for learning-based control," *IEEE Transactions on Automatic Control*, 2021.
- [33] B. Chachuat, B. Srinivasan, and D. Bonvin, "Adaptation strategies for real-time optimization," *Computers & Chemical Engineering*, vol. 33, no. 10, pp. 1557–1567, 2009.
- [34] T. Raissi, N. Ramdani, and Y. Candau, "Set membership state and parameter estimation for systems described by nonlinear differential equations," *Automatica*, vol. 40, no. 10, pp. 1771–1777, 2004.
- [35] M. Karimshoushtari and C. Novara, "Design of experiments for nonlinear system identification: A set membership approach," *Automatica*, vol. 119, p. 109036, 2020.
- [36] J. H. Manton, P.-O. Amblard *et al.*, "A primer on reproducing kernel Hilbert spaces," *Foundations and Trends in Signal Processing*, vol. 8, no. 1–2, pp. 1–126, 2015.
- [37] M. Milanese and C. Novara, "Set membership identification of nonlinear systems," *Automatica*, vol. 40, no. 6, pp. 957–975, 2004.
- [38] G. E. Fasshauer, "Positive definite kernels: past, present and future," *Dolomites Research Notes on Approximation*, vol. 4, pp. 21–63, 2011.
- [39] Y. Zhang, J. Duchi, and M. Wainwright, "Divide and conquer kernel ridge regression," in *Conference on Learning Theory (COLT)*. PMLR, 2013, pp. 592–617.
- [40] D. R. Burt, C. E. Rasmussen, and M. van der Wilk, "Convergence of sparse variational inference in gaussian processes regression," *Journal of Machine Learning Research*, vol. 21, pp. 1–63, 2020.
- [41] S. Boyd, N. Parikh, and E. Chu, *Distributed optimization and statistical learning via the alternating direction method of multipliers*. Now Publishers Inc, 2011.
- [42] A. Beck, "On the convergence of alternating minimization for convex programming with applications to iteratively reweighted least squares and decomposition schemes," *SIAM Journal on Optimization*, vol. 25, no. 1, pp. 185–209, 2015.
- [43] C. E. Rasmussen and C. K. I. Williams, *Gaussian Processes for Machine Learning*. The MIT Press, 2006.
- [44] E. Parzen *et al.*, "An approach to time series analysis," *Annals of mathematical statistics*, vol. 32, no. 4, pp. 951–989, 1961.
- [45] 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.

TABLE I

AVERAGE DISTANCE BETWEEN THE UPPER AND LOWER BOUNDS FOR THE OPTIMAL (OPT) AND SUB-OPTIMAL (SUB) DETERMINISTIC CASES, AND THE GAUSSIAN PROCESS ALTERNATIVE (GP). MODERATE NOISE LEVEL (TRUE $\bar{\delta} = 1$), USING FACTORS OF 1, 1.5, AND 2 TO AUGMENT $\bar{\delta}$ AND Γ .

Γ		1200			1800			2400		
$\bar{\delta}$		1	1.5	2	1	1.5	2	1	1.5	2
Grid	opt	6.21	8.35	10.34	7.45	9.75	11.90	8.50	10.94	13.20
	sub	11.07	15.60	20.13	11.70	16.23	20.76	12.36	16.89	21.42
	gp	604.51	706.13	786.51	904.61	1055.89	1175.32	1204.71	1405.65	1564.12
Rand	opt	14.62	19.02	22.89	18.05	23.08	27.51	20.85	26.39	31.26
	sub	64.78	93.99	123.20	65.91	95.12	124.33	67.07	96.28	125.49
	gp	643.20	743.44	822.24	962.51	1111.67	1228.70	1281.82	1479.90	1635.17

TABLE II

AVERAGE DISTANCE BETWEEN THE UPPER AND LOWER BOUNDS FOR THE OPTIMAL (OPT) AND SUB-OPTIMAL (SUB) DETERMINISTIC CASES, AND THE GAUSSIAN PROCESS ALTERNATIVE (GP). HIGH NOISE LEVEL (TRUE $\bar{\delta} = 5$), USING FACTORS OF 1, 1.5, AND 2 TO AUGMENT $\bar{\delta}$ AND Γ .

Γ		1200			1800			2400		
$\bar{\delta}$		5	7.5	10	5	7.5	10	5	7.5	10
Grid	opt	20.29	28.57	36.39	22.54	31.31	39.58	24.41	33.56	42.17
	sub	49.15	71.79	94.44	49.81	72.46	95.11	50.48	73.13	95.78
	gp	1090.16	1247.34	1366.96	1624.19	1854.47	2028.24	2158.21	2461.60	2689.52
Rand	opt	39.95	53.43	65.41	47.00	62.15	75.57	52.76	69.32	83.89
	sub	312.44	458.51	604.57	313.61	459.68	605.74	314.79	460.85	606.91
	gp	1117.01	1268.40	1383.43	1664.18	1885.79	2052.67	2211.36	2503.18	2721.91

- [46] A. Lederer, J. Umlauf, and S. Hirche, "Uniform error bounds for gaussian process regression with application to safe control," in *Conference on Neural Information Processing Systems (NeurIPS)*, 2019.
- [47] F. Berkenkamp, M. Turchetta, A. Schoellig, and A. Krause, "Safe model-based reinforcement learning with stability guarantees," *Advances in neural information processing systems*, vol. 30, 2017.
- [48] C. Fiedler, C. W. Scherer, and S. Trimpe, "Practical and rigorous uncertainty bounds for gaussian process regression," in *Proceedings of the AAAI conference on artificial intelligence*, vol. 35, no. 8, 2021, pp. 7439–7447.
- [49] L. Hewing, J. Kabzan, and M. N. Zeilinger, "Cautious model predictive control using gaussian process regression," *IEEE Transactions on Control Systems Technology*, vol. 28, no. 6, pp. 2736–2743, 2019.
- [50] T. Koller, F. Berkenkamp, M. Turchetta, and A. Krause, "Learning-based model predictive control for safe exploration," in *2018 IEEE Conference on Decision and Control (CDC)*, 2018, pp. 6059–6066.
- [51] S. R. Chowdhury and A. Gopalan, "On kernelized multi-armed bandits," in *International Conference on Machine Learning*. PMLR, 2017, pp. 844–853.
- [52] 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*, vol. 58, no. 5, pp. 3250–3265, 2012.