

On Robustness of Kernel-Based Regularized System Identification

Mohammad Khosravi and Roy S. Smith

Automatic Control Laboratory, ETH Zürich
 {khosravm, rsmith}@control.ee.ethz.ch

Abstract

This paper presents a novel feature of the kernel-based system identification method. We prove that the regularized kernel-based approach for the estimation of a finite impulse response is equivalent to a robust least-squares problem with a particular uncertainty set defined in terms of the kernel matrix, and thus, it is called *kernel-based uncertainty set*. We provide a theoretical foundation for the robustness of the kernel-based approach to input disturbances. Based on robust and regularized least-squares methods, different formulations of system identification are considered, where the kernel-based uncertainty set is employed in some of them. We apply these methods to a case where the input measurements are subject to disturbances. Subsequently, we perform extensive numerical experiments and compare the results to examine the impact of utilizing kernel-based uncertainty sets in the identification procedure. The numerical experiments confirm that the robust least square identification approach with the kernel-based uncertainty set improves the robustness of the estimation to the input disturbances.

1 Introduction

The system identification problem, as initially introduced by [1], deals with building suitable mathematical models for dynamic systems based on measurement data. Due to the numerous important applications, the subject has received extensive attention [2] developing identification tools and techniques, mainly based on the theory of mathematical statistics and optimization. One of these ideas is employing a suitable regularization to improve the estimation quality by resolving the overfitting issue and also integrating prior knowledge in the estimated model [3–5]. The regularization plays this role by introducing a penalty for the model candidates not satisfying specific desired properties. For example, when the model is supposed to have low complexity, the system order is penalized by the rank or the nuclear norm of the Hankel matrix [6], or based on a similar argument, the atomic norm of a transfer function is utilized in [7].

In the seminal work by Pillonetto and De Nicolao [8], the idea of utilizing *kernel-based* regularization is introduced and led to a paradigm shift in system identification [9]. In this approach, the identification problem is framed as a regularized regression where the regularization term is defined based on the norm of a reproducing kernel Hilbert space (RKHS) [10] with a suitable structure. The kernel-based approach allows the employment of suitable Tikhonov-like regularizations and resolves the bias-variance trade-off issue which was not addressed appropriately in previous works [11]. In this framework, the model complexity tuning, the counterpart of model order selection in the classical approaches, is performed efficiently by estimating the hyperparameters determining the kernel and the regularization term [4]. Moreover, various types of prior knowledge and desired features of

the model can be included in the estimated model by employing suitable forms of the kernel and imposing appropriate constraints on the regression problem. The stability and the smoothness of the impulse response can be enforced by utilizing kernel functions such as stable splines [4, 12]. By designing kernels based on particular filters, frequency domain attributes such as the time constant and the resonant frequency can be included in the estimated model [13]. In [14, 15], the positivity of the system is addressed by imposing structural constraints in the estimation problem. The prior knowledge on the DC-gain of the system is considered in [16]. The kernel-based paradigm is extended to the systems with specific structures, e.g., for Hammerstein and Wiener systems [17, 18], networked systems [19], and periodic systems [20]. Identifying models with low complexity is addressed based on the idea of multi-kernel regularization and sparse hyperparameter selection [21, 22]. The idea of multiple regularizations together with advanced hyperparameter tuning techniques are used for improving the performance of model estimation [23–26]. Recently, it has been empirically observed that kernel-based approaches can improve the robustness of the estimation to the input disturbances [27].

Inspired by [28], we introduce a new aspect of kernel-based identification in this paper. More precisely, we show that the kernel-based regularized estimation of a finite impulse response is equivalent to a robust least-squares problem. This result provides a theoretical foundation for the robustness of the kernel-based approach with respect to the input disturbances. The uncertainty set obtained in the robust optimization problem has an interesting special shape and structure which is defined based on the regularization matrix. Accordingly, the set is called, the *kernel-based uncertainty set*. In order to study the nature of this set, we consider various identification approaches formulated in terms of robust and regularized least-squares, and then apply these methods to an estimation case in which the input measurements are subject to disturbances. By means of Monte Carlo experiments, we compare these approaches to examine the impact of employing kernel-based uncertainty sets. The numerical experiments show that a robust least-squares approach with a kernel-based uncertainty set, not only inherits the interesting features of the kernel-based identification approach, but also improves the robustness of the estimation with respect to the input disturbances.

2 Notation

The set of integers, the set of non-negative integer numbers, the set of real numbers, and the set of non-negative real numbers are denoted by \mathbb{Z} , \mathbb{Z}_+ , \mathbb{R} , and \mathbb{R}_+ , respectively. The n -dimensional Euclidean space is denoted by \mathbb{R}^n and the set of n by m matrices is $\mathbb{R}^{n \times m}$. The zero vector is denoted by $\mathbf{0}$ and the identity matrix is denoted by \mathbb{I} . For a matrix $A \in \mathbb{R}^n$, $\text{tr}(A)$ denotes the trace of A . For a non-singular matrix A , the transpose of A^{-1} is $A^{-\text{T}}$. Given a positive-definite matrix $K \in \mathbb{R}^{n \times n}$, we define an inner product on \mathbb{R}^n as $\langle a, b \rangle_K := a^{\text{T}} K b$, for any $a, b \in \mathbb{R}^n$. Subsequently, an induced norm, denoted by $\|\cdot\|_K$, is defined as $\|a\|_K := (a^{\text{T}} K a)^{\frac{1}{2}}$. When $K = \mathbb{I}$, we have the Euclidean norm which is denoted by $\|\cdot\|$. Similarly, we can define an inner product on $\mathbb{R}^{m \times n}$ denoted by $\langle A, B \rangle_K$ and defined as $\langle A, B \rangle_K := \text{tr}(A K B^{\text{T}})$, for any $A, B \in \mathbb{R}^{m \times n}$. The corresponding induced norm is shown by $\|\cdot\|_K$. When $K = \mathbb{I}$, we have the Frobenius norm denoted by $\|\cdot\|_{\text{F}}$. The expression $X \sim \mathcal{N}(\mu, \Sigma)$ says that random variable X has a Gaussian distribution with mean μ and covariance Σ . We denote the probability density and the conditional probability density by $p(\cdot)$ and $p(\cdot|\cdot)$, respectively.

3 Regularized System Identification

Let \mathcal{S} be a discrete-time single-input-single-output stable and causal linear time-invariant (LTI) system, with a transfer function $G_{\mathcal{S}}(q)$ defined as

$$G_{\mathcal{S}}(q) := \sum_{k=0}^{\infty} g_k q^k, \quad (1)$$

where $g_{\mathcal{S}} := (g_k)_{k \in \mathbb{Z}_+}$ denotes the impulse response of the system and q denotes the *forward shift operator*. Due to the stability of \mathcal{S} , we know that $g \in \ell_1(\mathbb{Z}_+)$, i.e.,

$$\sum_{k=0}^{\infty} |g_k| < \infty. \quad (2)$$

Therefore, for any $\epsilon > 0$, there exists $n_{\mathcal{S}}(\epsilon) \in \mathbb{Z}_+$ such that $\sum_{k=n_{\mathcal{S}}(\epsilon)}^{\infty} |g_k| < \epsilon$. Accordingly, the infinite impulse response (IIR) can be truncated at a sufficiently high order $n_g \in \mathbb{Z}_+$ to approximate the system with a finite-length impulse response (FIR) denoted by g and defined as

$$g := [g_0 \quad g_1 \quad g_2 \quad \dots \quad g_{n_g-1}]^T \in \mathbb{R}^{n_g}. \quad (3)$$

Thus, the finite impulse response (FIR) model of the system is as

$$G(q) := \sum_{k=0}^{n_g-1} g_k q^k. \quad (4)$$

3.1 Identification Problem: FIR Formulation

Let $u := (u_t)_{t \in \mathbb{Z}}$ be a bounded input signal given to the system \mathcal{S} . Also, let $y := (y_t)_{t \in \mathbb{Z}}$ be the corresponding output signal which is subject to measurement uncertainty. Accordingly, we have

$$y_t = \sum_{k=0}^{n_g-1} g_k u_{t-k} + w_t, \quad \forall t \in \mathbb{Z}, \quad (5)$$

where w_t corresponds to the measurement noise and the unmodeled part of the system. Note that when the output measurement is subject to bounded uncertainty, then due to the stability of the system y , is a bounded signal. For time instants $t = 0, 1, \dots, n_{\mathcal{D}} - 1$, let assume the inputs and the measured outputs of the system are given. We define the set of data, denoted by \mathcal{D} , as

$$\mathcal{D} := \left\{ (u_t, y_t) \mid t = 0, 1, \dots, n_{\mathcal{D}} - 1 \right\}. \quad (6)$$

Following this, the identification problem is defined as estimating the FIR model of the system, g , using \mathcal{D} .

3.2 Prediction Error Method

Let vector φ_t be defined as

$$\varphi_t := [u_t \quad u_{t-1} \quad \dots \quad u_{t-n_g+1}]^T \in \mathbb{R}^{n_g}, \quad (7)$$

for any $t \in \mathbb{Z}$. For a candidate FIR model g , one can introduce the one-step ahead prediction rule for the output at time instant t as $\hat{y}(t|g) := \varphi_t^T g$. Considering the set of data \mathcal{D} , the quality of

the prediction rule can be assessed by comparing the actual measured outputs with the predicted values. To this end, one can form an empirical loss, $\mathcal{V}_{\mathcal{D}} : \mathbb{R}^{n_{\mathcal{g}}} \rightarrow \mathbb{R}$, for evaluating the prediction quality over set of data \mathcal{D} , e.g., $\mathcal{V}_{\mathcal{D}}$ can be defined as sum of squared errors of predictions, i.e., one has

$$\mathcal{V}_{\mathcal{D}}(\mathbf{g}) := \sum_{t=0}^{n_{\mathcal{D}}-1} (y_t - \hat{y}(t|\mathbf{g}))^2. \quad (8)$$

Then, one can estimate the FIR model by minimizing the empirical loss $\mathcal{V}_{\mathcal{D}}$. Define vectors \mathbf{y} and \mathbf{w} respectively as

$$\mathbf{y} := [y_0 \quad \cdots \quad y_{n_{\mathcal{D}}-1}]^{\top} \in \mathbb{R}^{n_{\mathcal{D}}}, \quad (9)$$

and

$$\mathbf{w} := [w_0 \quad \cdots \quad w_{n_{\mathcal{D}}-1}]^{\top} \in \mathbb{R}^{n_{\mathcal{D}}}. \quad (10)$$

Then, due to (5), one can easily see that

$$\mathbf{y} = \Phi \mathbf{g} + \mathbf{w}, \quad (11)$$

where $\Phi \in \mathbb{R}^{n_{\mathcal{D}} \times n_{\mathcal{g}}}$ is a Toeplitz matrix defined as follows

$$\Phi := \begin{bmatrix} \varphi_0^{\top} \\ \varphi_1^{\top} \\ \vdots \\ \varphi_{n_{\mathcal{D}}-1}^{\top} \end{bmatrix} = \begin{bmatrix} u_0 & u_{-1} & \cdots & u_{-n_{\mathcal{g}}+1} \\ u_1 & u_0 & \cdots & u_{-n_{\mathcal{g}}+2} \\ \vdots & \vdots & \ddots & \vdots \\ u_{n_{\mathcal{D}}-1} & u_{n_{\mathcal{D}}-2} & \cdots & u_{n_{\mathcal{D}}-n_{\mathcal{g}}} \end{bmatrix}. \quad (12)$$

Consequently, regarding the empirical loss, we have $\mathcal{V}_{\mathcal{D}}(\mathbf{g}) = \|\mathbf{y} - \Phi \mathbf{g}\|^2$. Here, for the sake of simplicity, one can assume that the system is initially at rest, i.e., $u_t = 0$, for all $t < 0$, and thus, given set of data \mathcal{D} , matrix Φ is known. Accordingly, estimating \mathbf{g} by minimizing the empirical loss is well defined and leads to the following least-squares (LS)

$$\begin{aligned} \mathbf{g}^{\text{LS}} &:= \operatorname{argmin}_{\mathbf{g} \in \mathbb{R}^{n_{\mathcal{g}}}} \|\mathbf{y} - \Phi \mathbf{g}\|^2 \\ &= (\Phi^{\top} \Phi)^{-1} \Phi^{\top} \mathbf{y}, \end{aligned} \quad (13)$$

where the last equality holds when $\Phi^{\top} \Phi$ is non-singular. When $\Phi^{\top} \Phi$ is not full-rank, the solution is not unique and the estimation problem is ill-posed. This issue happens for example when $n_{\mathcal{D}}$ is small. If the condition number of $\Phi^{\top} \Phi$ is high, then the estimation \mathbf{g}^{LS} is sensitive to noise. These issues can be resolved by including an appropriate regularization in (13) [4, 12].

3.3 Regularization Method

Let \mathbf{K} be a positive definite matrix. Define the regularization function $\mathcal{R} : \mathbb{R}^{n_{\mathcal{g}}} \rightarrow \mathbb{R}$ as $\mathcal{R}(\mathbf{g}) = \mathbf{g}^{\top} \mathbf{K}^{-1} \mathbf{g}$, for any $\mathbf{g} \in \mathbb{R}^{n_{\mathcal{g}}}$. In order to introduce the regularized estimation problem, we add this term to the empirical loss and obtain a regularized loss function $\mathcal{J} : \mathbb{R}^{n_{\mathcal{g}}} \rightarrow \mathbb{R}$ as

$$\mathcal{J}(\mathbf{g}) := \mathcal{V}_{\mathcal{D}}(\mathbf{g}) + \lambda \mathcal{R}(\mathbf{g}), \quad \forall \mathbf{g} \in \mathbb{R}^{n_{\mathcal{g}}}, \quad (14)$$

where $\lambda > 0$ is the regularization weight.

Let $\mathcal{V}_{\mathcal{G}}$ be defined as in (8). Then, the regularized estimation cost in (14) can be written as

$$\mathcal{J}(\mathbf{g}) := \|\mathbf{y} - \Phi\mathbf{g}\|^2 + \lambda\mathbf{g}^\top\mathbf{K}^{-1}\mathbf{g}. \quad (15)$$

Then, the regularized estimation of the FIR is

$$\begin{aligned} \mathbf{g}^{\text{Reg}} &:= \underset{\mathbf{g} \in \mathbb{R}^{n_{\mathcal{G}}}}{\text{argmin}} \|\mathbf{y} - \Phi\mathbf{g}\|^2 + \lambda\mathbf{g}^\top\mathbf{K}^{-1}\mathbf{g}, \\ &= (\Phi^\top\Phi + \lambda\mathbf{K}^{-1})\Phi^\top\mathbf{y}. \end{aligned} \quad (16)$$

The regularization matrix \mathbf{K} , also known as the *kernel matrix*, integrates in to the estimated impulse response additionally available prior information and desired attributes such as stability and smoothness. Additionally, it can reduce the variance of the estimation and resolve the potential ill-posedness of the problem and the bias-variance trade-off issue. In the literature, various methods are introduced for determining the regularization matrix \mathbf{K} [4, 12, 13, 29], e.g., by employing common kernels such as *tuned/correlated* (TC) and *stable spline* (SS) [4], or designing \mathbf{K} can be based on multiple regularizations [21–24, 26].

4 Robustness of Kernel-Based Regularized System Identification

In this section, we show that the kernel-based regularized identification is equivalent to a robust least-squares estimation with a specific structure.

Theorem 1. *There exist $\rho > 0$ such that the regularized estimation of the impulse response derived in (16) is the solution of the following robust least-squares problem*

$$\min_{\mathbf{g} \in \mathbb{R}^{n_{\mathcal{G}}}} \max_{\substack{\Delta \in \mathbb{R}^{n_{\mathcal{G}} \times n_{\mathcal{G}}} \\ \|\Delta\|_{\mathbf{K}} \leq \rho}} \|(\Phi + \Delta)\mathbf{g} - \mathbf{y}\|. \quad (17)$$

Proof. Let \mathbf{R} be the square root of \mathbf{K} or the lower triangular matrix in the Cholesky decomposition of \mathbf{K} . Then, the inner problem in (17) can be shown as to be

$$\max_{\substack{\Delta \in \mathbb{R}^{n_{\mathcal{G}} \times n_{\mathcal{G}}} \\ \|\Delta\mathbf{R}\|_{\mathbf{F}} \leq \rho}} \|(\Phi + \Delta)\mathbf{g} - \mathbf{y}\|^2. \quad (18)$$

Due to Lemma 3 (see Appendix A.2), the optimal value of (18) equals $\|\Phi\mathbf{g} - \mathbf{y}\| + \|\rho\mathbf{R}^{-1}\mathbf{x}\|$. Accordingly, the inner problem is equivalent to the following convex program

$$\begin{aligned} \min_{a,b} \quad & a \\ \text{s.t.} \quad & \|\Phi\mathbf{g} - \mathbf{y}\| \leq a - b, \\ & \|\rho\mathbf{R}^{-1}\mathbf{g}\| \leq b. \end{aligned} \quad (19)$$

Consequently, the robust optimization problem (17) can be written as in the following equivalent form

$$\begin{aligned} \min_{\substack{\mathbf{g} \in \mathbb{R}^{n_{\mathcal{G}}} \\ a,b \in \mathbb{R}}} \quad & a \\ \text{s.t.} \quad & \|\Phi\mathbf{g} - \mathbf{y}\| \leq a - b, \\ & \|\rho\mathbf{R}^{-1}\mathbf{g}\| \leq b. \end{aligned} \quad (20)$$

Note that at the optimal solution (g^*, a^*, b^*) , we have $a^* = \|\Phi g^* - y\| + \rho\|\mathbf{R}^{-1}g^*\|$, where a^* is the optimal value of (20) as well. Defining the vector x as $x := [g^T \quad a \quad b]^T$, we can re-write the problem (20) in the following form

$$\begin{aligned} \min_{x \in \mathbb{R}^{n_g+2}} \quad & [\mathbf{0}^T \quad 1 \quad 0] x \\ \text{s.t.} \quad & \|\Phi \mathbf{0} \mathbf{0}\| x - y\| \leq [\mathbf{0}^T \quad 1 \quad -1] x, \\ & \|\rho\mathbf{R}^{-1} \mathbf{0} \mathbf{0}\| x\| \leq [\mathbf{0}^T \quad 0 \quad 1] x. \end{aligned} \quad (21)$$

If $\epsilon > 0$, then $x = [\mathbf{0}^T \quad \|y\| + 1 + \epsilon \quad 1]^T$ is strictly feasible for (21), and also, the problem is bounded. Therefore, the Slater's conditions hold. The dual of (21) is as follows

$$\begin{aligned} \max_{s,t,z,w} \quad & y^T z \\ \text{s.t.} \quad & \begin{bmatrix} \Phi^T \\ \mathbf{0}^T \\ \mathbf{0}^T \end{bmatrix} z + \begin{bmatrix} \rho\mathbf{R}^{-T} \\ \mathbf{0}^T \\ \mathbf{0}^T \end{bmatrix} w + \begin{bmatrix} \mathbf{0} \\ 1 \\ -1 \end{bmatrix} s + \begin{bmatrix} \mathbf{0} \\ 0 \\ 1 \end{bmatrix} t = \begin{bmatrix} \mathbf{0} \\ 1 \\ 0 \end{bmatrix}, \\ & \|z\| \leq s, \\ & \|w\| \leq t, \\ & z \in \mathbb{R}^{n_\mathcal{D}}, w \in \mathbb{R}^{n_g}, s, t \in \mathbb{R}, \end{aligned}$$

which simplifies to the following optimization problem

$$\begin{aligned} \max_{\substack{z \in \mathbb{R}^{n_\mathcal{D}} \\ w \in \mathbb{R}^{n_g}}} \quad & y^T z \\ \text{s.t.} \quad & \Phi^T z + \rho\mathbf{R}^{-T} w = \mathbf{0}, \\ & \|z\| \leq 1, \\ & \|w\| \leq 1. \end{aligned} \quad (22)$$

Note that (22) is feasible and bounded. Let (z^*, w^*) be the optimal solution for the dual problem (22). Due to strong duality, we know that $a^* = y^T z^*$. Accordingly, for the optimal solutions of (21) and (22), we have that

$$\begin{aligned} \|\Phi g^* - y\| + \rho\|\mathbf{R}^{-1}g^*\| &= y^T z^* \\ &= (y - \Phi g^*)^T z^* + g^{*\top} \Phi^T z^* \\ &= (y - \Phi g^*)^T z^* - \rho g^{*\top} \mathbf{R}^{-T} w^*. \end{aligned} \quad (23)$$

Since $\|z^*\|, \|w^*\| \leq 1$, due to Cauchy-Schwartz inequality, (23) holds if and only if

$$z^* = \frac{y - \Phi g^*}{\|y - \Phi g^*\|}, \quad w^* = -\frac{\mathbf{R}^{-1}g^*}{\|\mathbf{R}^{-1}g^*\|}. \quad (24)$$

Following this and due to the equality constraint in (22), one can see that

$$\begin{aligned} \mathbf{0} &= \Phi^T z^* + \rho\mathbf{R}^{-T} w^* \\ &= \Phi^T \frac{y - \Phi g^*}{\|y - \Phi g^*\|} - \rho\mathbf{R}^{-T} \frac{\mathbf{R}^{-1}g^*}{\|\mathbf{R}^{-1}g^*\|} \\ &= \frac{\Phi^T y - \Phi^T \Phi g^*}{\|y - \Phi g^*\|} - \frac{\rho\mathbf{K}^{-1}g^*}{\|\mathbf{R}^{-1}g^*\|}. \end{aligned} \quad (25)$$

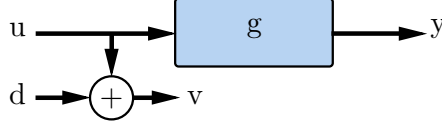


Figure 1: System with disturbance in input measurements.

Rearranging the terms in (25), one has

$$\Phi^\top \Phi g^* + \frac{\rho \|y - \Phi g^*\|}{\|R^{-1} g^*\|} K^{-1} g^* = \Phi^\top y. \quad (26)$$

Subsequently, it follows that

$$g^* = \left(\Phi^\top \Phi + \mu K^{-1} \right)^{-1} \Phi^\top y. \quad (27)$$

where $\mu = \frac{\rho \|y - \Phi g^*\|}{\|R^{-1} g^*\|}$. By choosing ρ appropriately, one can set $\mu = \lambda$. Consequently, due to (16), we have $g^* = g^{\text{Reg}}$. This concludes the proof. \blacksquare

Remark 1. From the proof of Theorem 1, one can see that

$$\rho := \lambda \frac{\|R^{-1} g^{\text{Reg}}\|}{\|\Phi g^{\text{Reg}} - y\|} = \lambda \frac{\left(g^{\text{Reg}\top} K^{-1} g^{\text{Reg}} \right)^{\frac{1}{2}}}{\|\Phi g^{\text{Reg}} - y\|}. \quad (28)$$

Accordingly, when the model fits well to the data, the residuals as well as $\|\Phi g^{\text{Reg}} - y\|$ are small. Subsequently ρ has a large value and the estimation is robust.

Remark 2. The main ingredient in the definition of the optimization problem (17) is the kernel matrix K . Accordingly, we call this method kernel-based robust least-squares.

5 Kernel-Based Uncertainty Set

In this section, we further study the *kernel-based uncertainty set* introduced in Section 4.

Let \mathcal{U}_ρ denote the uncertainty set employed in (17), i.e., \mathcal{U}_ρ is defined as

$$\mathcal{U}_\rho := \left\{ \Delta \in \mathbb{R}^{n_\varphi \times n_g} \mid \|\Delta\|_K^2 = \text{tr}(\Delta K \Delta^\top) \leq \rho^2 \right\}. \quad (29)$$

One can see that \mathcal{U}_ρ is a hyperball in the inner product space $(\mathbb{R}^{n_\varphi \times n_g}, \langle \cdot, \cdot \rangle_K)$ which is centered at the origin and has radius ρ . The inner product $\langle \cdot, \cdot \rangle_K$ is defined based on the kernel matrix K and determines the geometry of the space as well as the shape and the structure of \mathcal{U}_ρ .

In order to investigate the nature of the uncertainty set \mathcal{U}_ρ , we consider the identification problem where the measurements of the input signal are subject to disturbances, as shown in Figure 1. More precisely, at time instant t , the measured input is $v_t = u_t + d_t$, where u_t is the value of true input and d_t is the value of measurement disturbance, and also, let $y_t = \sum_{k=0}^{n_g-1} g_k u_{t-k}$ be the noiseless measured output of the system, for any $t = 0, \dots, n_\varphi - 1$. Let $\Psi, \Delta \in \mathbb{R}^{n_g \times n_\varphi}$ be Toeplitz matrices defined similarly to Φ as in (12) based on the measured inputs, disturbances and the unknown history of the input signal. Therefore, we have $\Psi = \Phi + \Delta$ and $y = \Phi g$, and subsequently, one can write y in two equivalent forms,

$$y = (\Psi - \Delta)g, \quad (30)$$

and

$$y = \Psi g + (-\Delta g). \quad (31)$$

Motivated by (30) and (31), one may propose variations on robust and regularized least-squares methods for estimating g . In the following, we provide the details of some these formulations.

1) Ordinary Least-Squares (LS): Motivated by (31) and based on a discussion similar to Section 3.2, one may propose a least-squares approach to estimate g as following

$$g^{\text{LS}} := \underset{g \in \mathbb{R}^{n_g}}{\operatorname{argmin}} \|y - \Psi g\|^2. \quad (32)$$

2) Regularized Least-Squares (RegLS): Similar to the least-squares case and the arguments provided in Section 3.3, a regularized least-squares approach can be proposed for estimating g as following

$$g^{\text{Reg}} := \underset{g \in \mathbb{R}^{n_g}}{\operatorname{argmin}} \|y - \Psi g\|^2 + \lambda g^T K^{-1} g. \quad (33)$$

3) Robust Least-Squares (RLS): Considering (30) and (31) along with the least-squares discussion, one may propose an estimation approach formulated as a standard robust least-squares problem as following

$$g^{\text{RLS}} := \underset{g \in \mathbb{R}^{n_g}}{\operatorname{argmin}} \left[\max_{\substack{\Delta \in \mathbb{R}^{n_\varphi \times n_g} \\ \|\Delta\|_{\text{F}} \leq \rho}} \|y - (\Psi - \Delta)g\|^2 \right]. \quad (34)$$

4) Structured Robust Least-Squares (SRLS): Since the uncertainty matrix Δ in (30) has a Toeplitz structure, we may formulate a structured robust least-squares for the estimation of g . To this end, given $\rho \in \mathbb{R}_+$, we define the *structured uncertainty set*, \mathcal{S}_ρ , as

$$\mathcal{S}_\rho := \left\{ \sum_{k=-n_g+1}^{n_\varphi-1} \delta_k E^{(k)} \mid \sum_{k=-n_g+1}^{n_\varphi-1} \delta_k^2 \leq \rho \right\},$$

where for each $k = -n_g + 1, \dots, n_\varphi - 1$, $E^{(k)} \in \mathbb{R}^{n_\varphi \times n_g}$ is a Toeplitz matrix with entries in $\{0, 1\}$ such that $E_{i,j}^{(k)} = 1$ only when $i - j = k$. Then, the estimation problem is

$$g^{\text{SRLS}} := \underset{g \in \mathbb{R}^{n_g}}{\operatorname{argmin}} \left[\max_{\Delta \in \mathcal{S}_\rho} \|y - (\Psi - \Delta)g\|^2 \right]. \quad (35)$$

5) Kernel-based Robust Least-Squares (KRLS): According to (30) and the discussion in Section 4 introducing the kernel-based uncertainty set \mathcal{U}_ρ (see equation (29)), we can formulate an estimation approach for g as the following robust least-squares problem

$$g^{\text{KRLS}} := \underset{g \in \mathbb{R}^{n_g}}{\operatorname{argmin}} \left[\max_{\Delta \in \mathcal{U}_\rho} \|y - (\Psi - \Delta)g\|^2 \right]. \quad (36)$$

6) Robust Regularized Least-Squares (RRegLS): By including regularization in RLS method, we obtain a robust regularized least-squares estimation approach as follows,

$$g^{\text{RReg}} := \underset{g \in \mathbb{R}^{n_g}}{\operatorname{argmin}} \left[\max_{\substack{\Delta \in \mathbb{R}^{n_\varphi \times n_g} \\ \|\Delta\|_{\text{F}} \leq \rho}} \|y - (\Psi - \Delta)g\|^2 + \lambda g^T K^{-1} g \right]. \quad (37)$$

7) Structured Robust Regularized Least-Squares (SRRegLS): Similar to the previous case, by considering the regularization in SRLS approach, one can obtain a structured robust regularized least-squares estimation as the following optimization problem

$$g^{\text{SRReg}} := \underset{g \in \mathbb{R}^{n_g}}{\operatorname{argmin}} \left[\max_{\Delta \in \mathcal{S}_\rho} \|y - (\Psi - \Delta)g\|^2 + \lambda g^\top K^{-1}g \right]. \quad (38)$$

8) Kernel-based Robust Regularized Least-Squares (KRRegLS): One can add the regularization term in KRLS approach, and subsequently, obtain a kernel-based robust regularized least-squares method for the estimation of g as the following optimization problem

$$g^{\text{KRReg}} := \underset{g \in \mathbb{R}^{n_g}}{\operatorname{argmin}} \left[\max_{\Delta \in \mathcal{U}_\rho} \|y - (\Psi - \Delta)g\|^2 + \lambda g^\top K^{-1}g \right]. \quad (39)$$

Remark 3. *Each of the above optimization problems is either convex or can be reformulated as a convex program.*

The differences between these formulations depend on the choice of uncertainty set and also whether a kernel-based regularization is employed or not. In order to investigate the impact of the kernel-based uncertainty set, employed in KRLS and KRRegLS, we compare the performance of these approaches by means of a Monte Carlo numerical experiment. To this end, we consider the following system [9],

$$G(q) = \frac{0.02008 + 0.04017q^{-1} + 0.02008q^{-2}}{1 - 1.561q^{-1} + 0.6414q^{-2}}, \quad (40)$$

and perform 1000 experiment runs. In each of these experiments, the system is simulated with a PRBS signal of length $n_\varphi = 127$ and the input is measured with measurement disturbance $d_t \sim \mathcal{N}(0, \sigma_d^2)$ where $\sigma_d = 0.1$, for $t = 0, \dots, n_\varphi - 1$. The output measurement is taken to be noiseless. In order to identify the system, we approximate G with a FIR g of length $n_g = 80$, and then apply each of the estimation methods formulated above. The value of ρ is calculated based on the true disturbances.

Figure 2 shows the histograms of normalized root mean squared errors (RMSE), $\frac{\|\hat{g} - g\|}{\|g\|}$, for each of the above methods. For all of the histograms, the x-axis range is taken as $[0, 0.6]$ to visualize the results better, however, the RLS and RRegLS have values above 0.6. In order to evaluate the quality of estimation, we employ R-squared metric defined as follows

$$R^2(\hat{g}) := 100 \times \left[1 - \frac{\|\hat{g} - g\|}{\|g - \bar{g}\|} \right], \quad (41)$$

where $\bar{g} \in \mathbb{R}^{n_g}$ is a vector such that each of its entries equals to $\frac{1}{n_g} \sum_{k=0}^{n_g-1} g_k$, i.e., the average of entries in g . Figure 3 demonstrates and compares the fitting results. The values of bias, variance and MSE for each of the above estimation methods are provided in Table 1.

6 Discussion

Based on the results of Monte Carlo experiments, we have the following observations:

- From Figure 2, we observe similar behaviors for RegLS and KRLS, which is expected due to Theorem 1. Moreover, the resulting variances are considerably smaller comparing to LS case demonstrating the robustness of the RegLS and KRLS approaches.

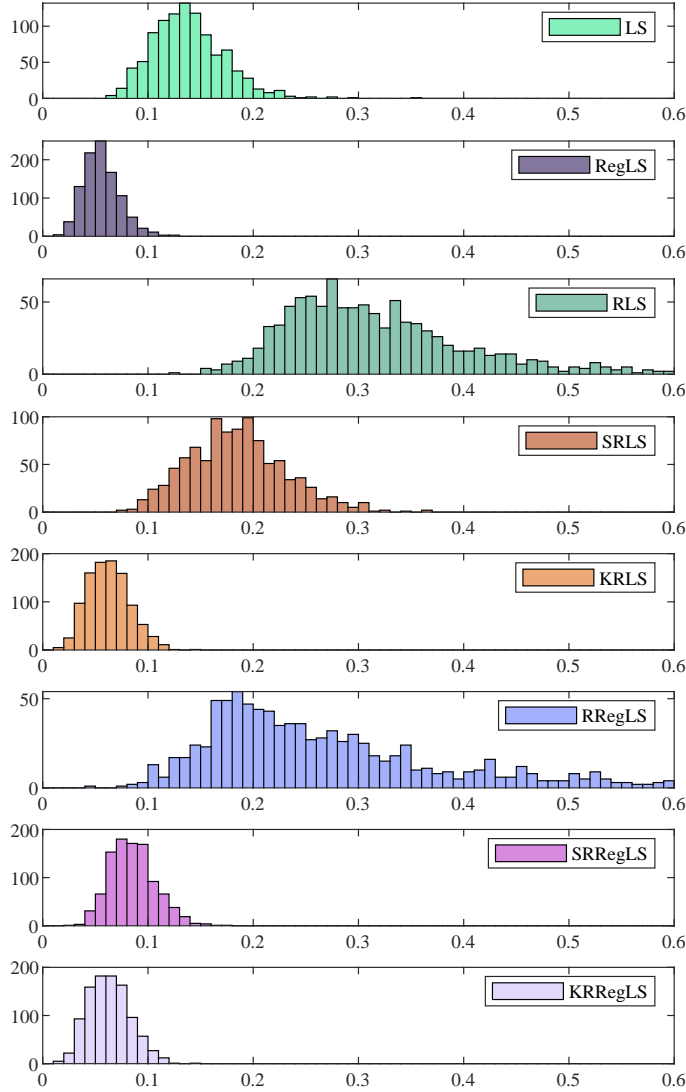


Figure 2: Histograms of normalized RMSE for different estimation methods.

- Note from Figure 2, Figure 3 and Table 1 that naïve formulations of robust optimization do not necessarily improve the robustness with respect to the input measurement disturbance. Indeed, the performance of RLS and RRegLS is significantly worse than LS. This can be explained by the nature of robust optimization which attempts to minimize the worst case cost. Accordingly, the shape and the structure of uncertainty set plays a major role in the performance of estimation approaches formulated as a robust optimization. This argument is supported by the performance improvements when RLS is compared to SRLS, or RRegLS is compared to SRRegLS.
- Comparing the estimation performances of RLS, SRLS and KRLS, one may conjecture that the right choices of uncertainty set for the estimation problems concerning system identification should have the same structure as the kernel-based uncertainty set. Comparing the performance of RRegLS, SRRegRS, and KRRegLS provides further support for this argument.

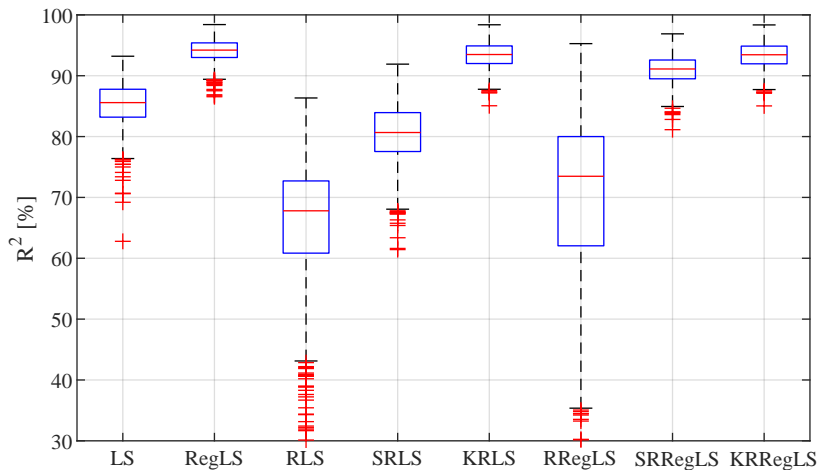


Figure 3: Comparison of fitting performances.

	Bias [$\times 10^4$]	Var [$\times 10^4$]	MSE [$\times 10^4$]	Unc. Set	Reg.
LS	0.56	22.0	22.3	-	-
RegLS	0.31	3.67	3.77	-	✓
RLS	10.66	23.7	137.4	standard	-
SRLS	4.16	22.5	39.8	structured	-
KRLS	1.68	1.91	4.74	kernel-based	-
RRegLS	10.35	47.3	154.4	standard	✓
SRRegLS	1.76	5.33	8.52	structured	✓
KRRegLS	1.70	1.90	4.79	kernel-based	✓

Table 1: Statistics for different methods.

- By comparing RLS and RRegLS, one can see that the choice of uncertainty set plays a more dominant role than the regularization term in the estimation method.
- From Table 1, one can see that the KRLS and KRRegLS methods have the smallest estimation variance and subsequently, maximal robustness to the input measurement disturbances. This observation highlights the impact of kernel-based uncertainty set. Also, the bias-variance trade-off for these approaches is close to the RegLS method. Moreover, one can see that including regularization in KRLS which results in KRRegLS does not improve the estimation performance significantly, at least for the current setting where the output measurement is noiseless.

7 Conclusions

In this paper, we have shown a novel feature of the kernel-based system identification method concerning robustness to the input disturbances. We have proved that the regularized kernel-based approach can be reformulated as a robust least-squares problem with an uncertainty set defined based on the kernel-matrix. Using extensive numerical experiments and comparisons, we have studied the nature of this new uncertainty set. It has been verified that the robust least

square identification approach with the kernel-based uncertainty set is robust with respect to input disturbances and retains other features of the kernel-based approach as well.

A Appendix

A.1 Duality for Second-Order Cone Programming

Consider the following second-order cone programming

$$\begin{aligned} \min_{\mathbf{x} \in \mathbb{R}^n} \quad & \mathbf{c}^\top \mathbf{x} \\ \text{s.t.} \quad & \|\mathbf{A}\mathbf{x} - \mathbf{d}\| \leq \mathbf{a}^\top \mathbf{x}, \\ & \|\mathbf{B}\mathbf{x}\| \leq \mathbf{b}^\top \mathbf{x}, \end{aligned} \tag{P}$$

where n, m, k are positive integers, $\mathbf{a}, \mathbf{b}, \mathbf{c} \in \mathbb{R}^n$, $\mathbf{A} \in \mathbb{R}^{m \times n}$, $\mathbf{B} \in \mathbb{R}^{k \times n}$, and $\mathbf{d} \in \mathbb{R}^m$. Then, the dual of (P) is the following convex program

$$\begin{aligned} \max_{s, t, \mathbf{z}, \mathbf{w}} \quad & \mathbf{d}^\top \mathbf{z} \\ \text{s.t.} \quad & \mathbf{A}^\top \mathbf{z} + \mathbf{B}^\top \mathbf{w} + \mathbf{a} s + \mathbf{b} t = \mathbf{c}, \\ & \|\mathbf{z}\| \leq s, \\ & \|\mathbf{w}\| \leq t, \\ & \mathbf{z} \in \mathbb{R}^m, \mathbf{w} \in \mathbb{R}^k, s, t \in \mathbb{R}. \end{aligned} \tag{D}$$

Note that when the dual problem (D) is feasible and bounded, then the primal problem (P) is also feasible and bounded. Moreover, they attain same optimal values when strong duality holds. For more details, see [30].

A.2 Supplementary Lemmas

Lemma 2. *Let $\mathbf{a} \in \mathbb{R}^n$ and $\mathbf{D} \in \mathbb{R}^{m \times n}$. Then, we have $\|\mathbf{D}\mathbf{a}\|^2 \leq \|\mathbf{D}\|_{\mathbb{F}}^2 \|\mathbf{a}\|^2$. The equality occurs iff there exists vector $\mathbf{c} \in \mathbb{R}^n$ such that $\mathbf{D} = \mathbf{c} \mathbf{a}^\top$.*

Proof. We know that

$$\begin{aligned} \|\mathbf{D}\mathbf{a}\|^2 &= \sum_{i=1}^m \left(\sum_{j=1}^n d_{ij} a_j \right)^2 \\ &\leq \sum_{i=1}^m \left(\sum_{j=1}^n d_{ij}^2 \right) \|\mathbf{a}\|^2 = \|\mathbf{D}\|_{\mathbb{F}}^2 \|\mathbf{a}\|^2, \end{aligned} \tag{42}$$

where the Cauchy-Schwartz inequality is used in the second step. For any $i = 1, \dots, m$, due to the equality condition for Cauchy-Schwartz theorem, we know that $(\sum_{j=1}^n d_{ij} a_j)^2 = (\sum_{j=1}^n d_{ij}^2) \|\mathbf{a}\|^2$ iff there exists $c_i \in \mathbb{R}$ such that we have $[d_{i1} \ \dots \ d_{in}] = c_i \mathbf{a}^\top$. Therefore, the equality occurs in (42) iff there exist $\mathbf{c} \in \mathbb{R}^n$ such that $\mathbf{D} = \mathbf{c} \mathbf{a}^\top$. ■

Lemma 3. *Let $\rho > 0$, $\mathbf{a} \in \mathbb{R}^n$, $\mathbf{b} \in \mathbb{R}^m$ and $\mathbf{R} \in \mathbb{R}^{n \times n}$ be an invertible matrix. Then, for the following optimization*

$$\max_{\substack{\Delta \in \mathbb{R}^{m \times n} \\ \|\Delta \mathbf{R}\|_{\mathbb{F}} \leq \rho}} \|\Delta \mathbf{a} + \mathbf{b}\|, \tag{43}$$

the optimal value equals $\|\mathbf{b}\| + \|\rho \mathbf{R}^{-1} \mathbf{a}\|$.

Proof. Since for $\mathbf{a} = \mathbf{0}$, the result is straightforward, we assume $\mathbf{a} \neq \mathbf{0}$. Also, let $\mathbf{b} \neq \mathbf{0}$. We know that

$$\|\Delta\mathbf{a} + \mathbf{b}\| \leq \|\Delta\mathbf{a}\| + \|\mathbf{b}\|, \quad (44)$$

where the equality holds if there exists $\eta > 0$ such that $\mathbf{b} = \eta\Delta\mathbf{a}$. Moreover, from Lemma 2, we have that

$$\begin{aligned} \|\Delta\mathbf{a}\| &= \|\Delta\mathbf{R}\mathbf{R}^{-1}\mathbf{a}\| \\ &\leq \|\Delta\mathbf{R}\|_{\mathbb{F}}\|\mathbf{R}^{-1}\mathbf{a}\| \\ &\leq \rho\|\mathbf{R}^{-1}\mathbf{a}\|, \end{aligned} \quad (45)$$

for any $\Delta \in \mathbb{R}^{m \times n}$ such that $\|\Delta\mathbf{R}\|_{\mathbb{F}} \leq \rho$. The equality occurs in (45) iff there exists $\mathbf{c} \in \mathbb{R}^m$ such that $\Delta\mathbf{R} = \mathbf{c}(\mathbf{R}^{-1}\mathbf{a})^{\top}$ and $\|\Delta\mathbf{R}\|_{\mathbb{F}} = \rho$. From (44) and (45), we have

$$\|\Delta\mathbf{a} + \mathbf{b}\| \leq \left(\rho\|\mathbf{R}^{-1}\mathbf{a}\| + \|\mathbf{b}\| \right)^2, \quad (46)$$

for any $\Delta \in \mathbb{R}^{m \times n}$ such that $\|\Delta\mathbf{R}\|_{\mathbb{F}} \leq \rho$. Note that the right-hand side in (46) does not depend on Δ and therefore, it is an upper bound for (43). If the equality holds in (46), for a given Δ^* , the equality conditions mentioned above should be satisfied. More precisely, we need to have $\Delta^* = \mathbf{c}(\mathbf{R}^{-1}\mathbf{a})^{\top}\mathbf{R}^{-1}$, for some $\mathbf{c} \in \mathbb{R}^m$ such that $\|\mathbf{c}\|\|\mathbf{R}^{-1}\mathbf{a}\| = \rho$ and there exists $\eta > 0$ such that

$$\mathbf{b} = \eta\mathbf{c}(\mathbf{R}^{-1}\mathbf{a})^{\top}\mathbf{R}^{-1}\mathbf{a} = \eta\|\mathbf{R}^{-1}\mathbf{a}\|^2\mathbf{c}. \quad (47)$$

Consequently, these equality holds iff

$$\mathbf{c} = \frac{\rho\mathbf{b}}{\|\mathbf{b}\|\|\mathbf{R}^{-1}\mathbf{a}\|}, \quad \eta = \frac{\rho\|\mathbf{b}\|}{\|\mathbf{R}^{-1}\mathbf{a}\|}. \quad (48)$$

Therefore, (43) has a unique solution Δ^* defined as

$$\Delta^* := \frac{\rho\mathbf{b}(\mathbf{R}^{-1}\mathbf{a})^{\top}\mathbf{R}^{-1}}{\|\mathbf{b}\|\|\mathbf{R}^{-1}\mathbf{a}\|}. \quad (49)$$

Moreover, the optimal value of (43) equals to the right-hand side of (46). For the case $\mathbf{b} = \mathbf{0}$, due to inequality (45), we know that

$$\|\Delta\mathbf{a} + \mathbf{b}\| = \|\Delta\mathbf{a}\| \leq \rho\|\mathbf{R}^{-1}\mathbf{a}\| = \rho\|\mathbf{R}^{-1}\mathbf{a}\| + \|\mathbf{b}\|, \quad (50)$$

for any $\Delta \in \mathbb{R}^{m \times n}$ such that $\|\Delta\mathbf{R}\|_{\mathbb{F}} \leq \rho$. The equality holds for Δ^* iff there exists $\mathbf{c} \in \mathbb{R}^m$ such that $\Delta^*\mathbf{R} = \mathbf{c}(\mathbf{R}^{-1}\mathbf{a})^{\top}$ and $\|\Delta^*\mathbf{R}\|_{\mathbb{F}} = \rho$, or equivalently $\Delta^* = \mathbf{c}(\mathbf{R}^{-1}\mathbf{a})^{\top}\mathbf{R}^{-1}$ and $\|\mathbf{c}\| = \frac{\rho}{\|\mathbf{R}^{-1}\mathbf{a}\|}$. Replacing \mathbf{c} with $\frac{\rho\mathbf{u}}{\|\mathbf{R}^{-1}\mathbf{a}\|}$, where \mathbf{u} is a unit vector in \mathbb{R}^m , we obtain Δ^* as

$$\Delta^* := \frac{\rho\mathbf{u}(\mathbf{R}^{-1}\mathbf{a})^{\top}\mathbf{R}^{-1}}{\|\mathbf{R}^{-1}\mathbf{a}\|}. \quad (51)$$

This concludes the proof. ■

References

- [1] L. Zadeh, “On the identification problem,” *IRE Transactions on Circuit Theory*, vol. 3, no. 4, pp. 277–281, 1956.
- [2] L. Ljung, *System Identification: Theory for the User*. Prentice Hall, 1999.
- [3] M. Khosravi and R. S. Smith, “Convex nonparametric formulation for identification of gradient flows,” *IEEE Control Systems Letters*, vol. 5, no. 3, pp. 1097–1102, 2021.
- [4] 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.
- [5] M. Khosravi and R. S. Smith, “Nonlinear system identification with prior knowledge on the region of attraction,” *IEEE Control Systems Letters*, vol. 5, no. 3, pp. 1091–1096, 2021.
- [6] R. S. Smith, “Frequency domain subspace identification using nuclear norm minimization and Hankel matrix realizations,” *IEEE Transactions on Automatic Control*, vol. 59, no. 11, pp. 2886–2896, 2014.
- [7] P. Shah, B. N. Bhaskar, G. Tang, and B. Recht, “Linear system identification via atomic norm regularization,” in *Conference on Decision and Control*. IEEE, 2012, pp. 6265–6270.
- [8] G. Pillonetto and G. De Nicolao, “A new kernel-based approach for linear system identification,” *Automatica*, vol. 46, no. 1, pp. 81–93, 2010.
- [9] L. Ljung, T. Chen, and B. Mu, “A shift in paradigm for system identification,” *International Journal of Control*, vol. 93, no. 2, pp. 173–180, 2020.
- [10] N. Aronszajn, “Theory of reproducing kernels,” *Transactions of the American Mathematical Society*, vol. 68, no. 3, pp. 337–404, 1950.
- [11] A. Chiuso, “Regularization and Bayesian learning in dynamical systems: past, present and future,” *Annual Reviews in Control*, vol. 41, pp. 24–38, 2016.
- [12] T. Chen, “On kernel design for regularized LTI system identification,” *Automatica*, vol. 90, pp. 109–122, 2018.
- [13] A. Marconato, M. Schoukens, and J. Schoukens, “Filter-based regularisation for impulse response modelling,” *IET Control Theory & Applications*, vol. 11, pp. 194–204, 2016.
- [14] M. Zheng and Y. Ohta, “Positive FIR system identification using maximum entropy prior,” *IFAC-PapersOnLine*, vol. 51, no. 15, pp. 7–12, 2018.
- [15] M. Khosravi and R. S. Smith, “Kernel-based identification of positive systems,” in *Conference on Decision and Control*. IEEE, 2019, pp. 1740–1745.
- [16] Y. Fujimoto and T. Sugie, “Kernel-based impulse response estimation with a priori knowledge on the DC gain,” *IEEE Control Systems Letters*, vol. 2, no. 4, pp. 713–718, 2018.
- [17] R. S. Risuleo, G. Bottegal, and H. Hjalmarsson, “A nonparametric kernel-based approach to Hammerstein system identification,” *Automatica*, vol. 85, pp. 234–247, 2017.

- [18] R. S. Risuleo, F. Lindsten, and H. Hjalmarsson, “Bayesian nonparametric identification of Wiener systems,” *Automatica*, vol. 108, p. 108480, 2019.
- [19] K. R. Ramaswamy, G. Bottegal, and P. M. Van den Hof, “Local module identification in dynamic networks using regularized kernel-based methods,” in *Conference on Decision and Control*. IEEE, 2018, pp. 4713–4718.
- [20] M. Yin, A. Iannelli, M. Khosravi, A. Parsi, and R. S. Smith, “Linear time-periodic system identification with grouped atomic norm regularization,” *arXiv preprint arXiv:2003.06653*, 2020.
- [21] 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.
- [22] M. Khosravi, M. Yin, A. Iannelli, A. Parsi, and R. S. Smith, “Low-complexity identification by sparse hyperparameter estimation,” *IFAC-PapersOnLine*, 2020.
- [23] S. Hong, B. Mu, F. Yin, M. S. Andersen, and T. Chen, “Multiple kernel based regularized system identification with SURE hyper-parameter estimator,” *IFAC-PapersOnLine*, vol. 51, no. 15, pp. 13–18, 2018.
- [24] T. Chen, M. S. Andersen, B. Mu, F. Yin, L. Ljung, and S. J. Qin, “Regularized LTI system identification with multiple regularization matrix,” *IFAC-PapersOnLine*, vol. 51, no. 15, pp. 180–185, 2018.
- [25] B. Mu, T. Chen, and L. Ljung, “Asymptotic properties of generalized cross validation estimators for regularized system identification,” *IFAC-PapersOnLine*, vol. 51, no. 15, pp. 203–208, 2018.
- [26] M. Khosravi, A. Iannelli, M. Yin, A. Parsi, and R. S. Smith, “Regularized system identification: A hierarchical Bayesian approach,” *IFAC-PapersOnLine*, 2020.
- [27] T. Hiroe, K. Ide, R. Sase, and T. Sugie, “Kernel-based system identification improving robustness to input disturbances: A preliminary study on its application to AR models,” pp. 546–551, 2020.
- [28] L. El Ghaoui and H. Le Bret, “Robust solutions to least-squares problems with uncertain data,” *SIAM Journal on Matrix Analysis and Applications*, vol. 18, no. 4, pp. 1035–1064, 1997.
- [29] M. Zorzi and A. Chiuso, “The harmonic analysis of kernel functions,” *Automatica*, vol. 94, pp. 125–137, 2018.
- [30] S. Boyd and L. Vandenberghe, *Convex Optimization*. Cambridge university press, 2004.