# On kernel design for regularized LTI system identification [*]

Tianshi Chen

*School of Science and Engineering, The Chinese University of Hong Kong, Shenzhen, 518172, China, (e-mail: tschen@cuhk.edu.cn).*

**Abstract**

There are two key issues for the kernel-based regularization method: one is how to design a suitable kernel to embed in the kernel the prior knowledge of the LTI system to be identified, and the other one is how to tune the kernel such that the resulting regularized impulse response estimator can achieve a good bias-variance tradeoff. In this paper, we focus on the issue of kernel design. Depending on the type of the prior knowledge, we propose two methods to design kernels: one is from a machine learning perspective and the other one is from a system theory perspective. We also provide analysis results for both methods, which not only enhances our understanding for the existing kernels but also directs the design of new kernels.

*Key words:* System identification, regularization methods, kernel methods, kernel design, prior knowledge.

## 1 Introduction

Among diverse system identification problems, the linear time-invariant (LTI) system identification is a classical and fundamental problem. As well-known, an LTI system is uniquely characterized by its impulse response, and thus the LTI system identification is equivalent to its impulse response estimation that could be ill-conditioned in practice. So to tackle the LTI system identification, one needs to first find a way to make the ill-conditioned problem well-conditioned. Then one faces the core issue of system identification, i.e., how to construct a model estimator able to achieve a good bias-variance tradeoff, or equivalently, to suitably balance the adherence to the data and the model complexity.

There are different routes to handle these two issues. The

___

most widely used route is to adopt instead a parametric model structure with a finite number of parameters. The so-called maximum likelihood/prediction error method (ML/PEM) is one method along this route. It has optimal asymptotic properties and is the current standard method to LTI system identification. It first postulates a parametric model structure, e.g., a rational transfer function, then forms the prediction error criterion, and finally obtains the model estimator by minimizing the prediction error criterion; see e.g., [1,2]. The model complexity of the parametric model structure is governed by the number of parameters and tuned in a discrete manner. The bias-variance tradeoff issue is handled by cross validation or model structure selection criteria such as Akaike's information criterion (AIC), Bayesian information criterion (BIC), and etc, which correspond to combinatorial optimization problems; see e.g., [1, Chapter 16]. However, these techniques are not as reliable as expected when the data is short and/or has low signal-to-noise ratio; see e.g., [3,4].

Another route is to adopt regularization by adding an extra regularization term in the estimation criterion. The kernel-based regularization method (KRM) proposed in the seminal paper [5] and further developed in [6,3,7] is one method along this route. It first proposes a kernel parameterized by some hyper-parameters for the impulse response, then estimates these hyper-parameters with certain methods, and finally obtains the estimate of the impulse response by minimizing a regularized least squares criterion with the regularization term in quadratic form and the regularization matrix defined

through the kernel; see [4] for a recent survey. While the impulse response model is used, the underlying model structure is determined by the kernel. The model complexity is governed by the hyper-parameters of the kernel, and tuned in a continuous manner by e.g., the empirical Bayes method, Stein's unbiased risk estimator [4,8,9] and etc., which correspond to non-convex but non-combinatorial optimization problems.

The route to adopt regularization is by no means new, see e.g., [10], [11], [1, p. 504-505] and also [12] for a historic review, but no important progress along this route has been reported until [5]. The major obstacle is that it was unclear whether or not it is possible to design the regularization to embed the prior knowledge of the LTI system to be identified. The intriguing finding disclosed by the KRM is that when considering impulse response estimation problem, it is possible to design a regularization in quadratic form or equivalently a kernel to embed the prior knowledge of the *impulse response* of the LTI system to be identified. So far several kernels have been invented to embed various prior knowledge, e.g., [5,6,3,7,13,14,15,16,17,18,19] and some analysis results of the corresponding kernels have been derived, e.g., [5,3,7,13,14,15,16,17,18]. Still, there lack systematic ways to design kernels to embed in the kernel the prior knowledge of the impulse response of the LTI system to be identified.

In this paper, we try to develop systematic ways to design kernels. Clearly, this issue relates to the type of the prior knowledge. Interestingly, even for the same impulse response estimation problem, users with different background may come up with different types of prior knowledge because they may treat the impulse response in different ways. For instance, users from machine learning may treat the impulse response as a function, whose amplitude varies with a certain rate and decays exponentially to zero in the end; users from signals and systems may associate the impulse response with an LTI system that is stable and may be overdamped, underdamped, has multiple distinct time constants and resonant frequencies, and etc. Accordingly, we provide two methods to design kernels to embed the corresponding type of prior knowledge from a machine learning perspective and from a system theory perspective, respectively.

To develop the machine learning method, we first point out a common feature of the two most widely used single kernels, i.e., the stable spline (SS) kernel in [5] and the diagonal-correlated (DC) kernel in [3], that they belong to the class of the so-called amplitude modulated locally stationary (AMLS) kernel. The AMLS kernel is a multiplication of two kernels: a rank-1 kernel and a stationary kernel. It has a neat implication: the zero mean Gaussian process with the AMLS kernel as the covariance function can be seen as an amplitude modulated (by the rank-1 kernel) zero mean Gaussian processes with the stationary kernel as covariance function; see Section 3 for details. This finding leads to the machine learning method to design general kernels: design the rank-1 kernel and the stationary kernel to account for the decay and the varying rate of the impulse response, respectively.

To develop the system theory method, we recall the general guideline to design a kernel in [3] that is to let the kernel mimic the behavior of the optimal kernel [3, Thm. 1] and [4, Prop. 19]. Moreover, the prior knowledge for the impulse response, or equivalently, the LTI system to be identified, shall be made use of, since the optimal kernel depends on the true impulse response. Following this guideline and employing the multiplicative uncertainty configuration in robust control, see e.g. [20], we design the simulation induced (SI) kernel. In particular, the prior knowledge is embedded in the nominal model, the uncertainty is assumed to be stable and finally, the multiplicative uncertainty configuration is used to take into account both the nominal model and the uncertainty and is simulated with an impulsive input to get the SI kernel. Then we make analysis for SI kernels in terms of stability, maximum entropy property, and Markov property, with the SS and DC kernels as examples. In particular, we give conditions under which the SI kernel is stable, solves a suitably defined entropy maximization problem, and induces a Gaussian Markov process. The maximum entropy interpretation enhances our understanding about a kernel, and the Markov property of a kernel ensures that the inverse of the kernel matrix is banded and this special structure can be exploited to derive more stable and efficient implementation for the KRM, see e.g., [21,18].

The remaining part of this paper is organized as follows. In Section 2, the KRM is recapped. The machine learning and system theory methods to design kernels are introduced in Section 3 and 4, respectively. Then a case study is provided in Section 5 to demonstrate how to design kernels from the proposed two methods to model impulse responses with damped oscillation. Finally, this paper is concluded in Section 6 with a take-home message. The proofs of all propositions and theorems are given in the Appendix.

## 2 LTI System Identification with Kernel-based Regularization Method

### 2.1 Impulse Response Estimation

Consider a single-input-single-output, causal and LTI system described by

$$y(t) = (g * u)(t) + v(t), \quad t = T_s, 2T_s, \cdots, NT_s, \quad (1)$$

where $t$ is the time index, $y(t), u(t), v(t) \in \mathbb{R}$ and $g(t) \in \mathbb{R}$ are the measured output, input, disturbance, and impulse response of the LTI system at time $t$, respectively,

$(g*u)(t)$ is the convolution between the impulse response $g(\cdot)$ and the input $u(\cdot)$ evaluated at time $t$, and $T_s$ is the sampling period and for convenience chosen to be 1.

The system (1) represents both discrete-time (DT) system with $t = 1, 2, \cdots$, and continuous-time system (CT) with $t \geq 0$. In particular, we have

$$(g*u)(t) = \begin{cases} \sum_{\tau=1}^{\infty} g(\tau)u(t-\tau), & \text{DT}, \\ \int_0^{\infty} g(\tau)u(t-\tau)d\tau, & \text{CT}, \end{cases} \quad (2)$$

where the unmeasured portions of $u(t)$ are set to zero. Moreover, the system (1) is assumed to be stable, i.e., $g \in \ell^1$ for the DT case and $g \in \mathcal{L}^1$ for the CT case, where $\ell^1$ denotes the space of absolutely convergent sequences and $\mathcal{L}^1$ denotes the space of absolutely Lebesgue integrable functions on $t \geq 0$, and the disturbance $v(t)$ is assumed to be a white Gaussian noise with zero mean and variance $\sigma^2$ and independent of $u(t)$.

The goal of LTI system identification is to estimate the impulse response $g(t)$ as well as possible given the data $y(t), u(t)$ with $t = 1, 2, \cdots, N$ for the DT case and $y(t)$ with $t = 1, 2, \cdots, N$ and $u(t)$ with $t \geq 0$ for the CT case.

*2.2 Kernel-based Regularization Method*

As well-known [10], [1, p. 504-505], straightforward impulse response estimation, i.e.,

$$\underset{g}{\text{minimize}} \sum_{t=1}^{N} (y(t) - (g*u)(t))^2$$

could be an ill-conditioned problem in practice and one way to overcome this problem is to adopt regularization. Moreover, the recent progresses for KRM [5,6,3,7] show that if the regularization is well designed and tuned, the resulting regularized impulse response estimator can also achieve a good bias-variance tradeoff.

To introduce the KRM, we first recall the definition of the positive semidefinite kernel and its associated reproducing kernel Hilbert space (RKHS). Let $(X, d)$ be a metric space with $d$ being its metric. A function $k : X \times X \to \mathbb{R}$ is called a positive semidefinite kernel, if it is symmetric and satisfies $\sum_{i,j=1}^{m} a_i a_j k(x_i, x_j) \geq 0$ for any finite set of points $\{x_1, \cdots, x_m\} \subset X$ and $\{a_1, ..., a_m\} \subset \mathbb{R}$. As well-known from e.g., [22], to every positive semidefinite kernel $k$ there corresponds to one and only one class of functions with a unique determined inner product in it, leading to the so-called reproducing kernel Hilbert space (RKHS) $\mathcal{H}_k$ with $k$ as the reproducing kernel.

The KRM first introduces a suitable positive semidef-

inite kernel $k(t, s; \theta)$ [1] with $t, s \in X$, where for the DT case, $X = \mathbb{N}$ and for the CT case, $X = \{t | t \geq 0\}$, and $\theta \in \mathbb{R}^m$ is a hyper-parameter vector that contains the parameters used to parameterize the kernel, and then solves the following regularized least squares problem:

$$\hat{g}(t) = \underset{g \in \mathcal{H}_k}{\arg\min} \sum_{t=1}^{N} (y(t) - (g*u)(t))^2 + \gamma \|g\|_{\mathcal{H}_k}^2, \quad (3)$$

where $\|g\|_{\mathcal{H}_k}$ is the norm of $g$ in $\mathcal{H}_k$, and $\gamma > 0$ is a regularization parameter and controls the tradeoff between the data fit $\sum_{t=1}^{N}(y(t) - (g*u)(t))^2$ and the regularization term $\|g\|_{\mathcal{H}_k}^2$.

We will discuss the issue of kernels in the next subsections. For the time being, we assume that a suitable kernel $k(t, s; \theta)$ has been designed, but its hyper-parameter $\theta$ is left to be determined. The current most widely used method to determine $\theta$ is the so-called empirical Bayes method [23]. It first embeds the regularization in a Bayesian framework, and then estimates $\theta$ and possibly also the noise variance $\sigma^2$ by maximizing the marginal likelihood $p(Y|\eta)$, where $Y \in \mathbb{R}^N$ with $y(t)$ being its $t$th element, and $\eta$ could be $\theta$ or the concatenation of $\theta$ and $\sigma^2$. Specifically, we define $A(\theta) \in \mathbb{R}^{N \times N}$ with its $(t, s)$th element $A_{t,s}(\theta)$ defined by

$$a(t, s; \theta) = (k(t, \cdot; \theta) * u)(s), A_{t,s}(\theta) = (a(\cdot, s; \theta) * u)(t)$$

and moreover, we let $\Sigma(\eta) = A(\theta) + \sigma^2 I_N$, where $I_N$ is the $N$-dimensional identity matrix. Then we get

$$\underset{\eta \in \Gamma}{\text{minimize}} \, Y^T \Sigma(\eta)^{-1} Y + \log \det \Sigma(\eta),$$

where $\Gamma$ is a constraint set where we search for $\eta$. When an estimate of $\eta$ is obtained, the solution to (3) is given by the representer theorem [4, Theorem 3]:

$$\hat{g}(t) = \sum_{s=1}^{N} \hat{c}_s a(t, s; \theta), \quad (4)$$

where $\hat{c}_s$ is the $s$th element of $\hat{c} = (A(\theta) + \sigma^2 I_N)^{-1} Y$.

*2.3 Existing Single Kernels*

So far the two mostly widely used single kernels are the stable spline (SS) kernel [5] and the diagonal/correlated (DC) kernel [3]. The SS kernel is defined as

$$k^{\text{SS}}(t, s; \theta) = c\frac{1}{2} \exp(-\beta(t+s) - \beta \max\{t, s\})$$
$$- c\frac{1}{6} \exp(-3\beta \max\{t, s\}), \, \theta = [c \ \beta]^T, \, c, \beta \geq 0, \quad (5)$$

---

[1] Sometimes the dependence of $k(t, s; \theta)$ on $\theta$ is ignored.

and the DC kernel is defined as

$$k^{\mathrm{DC}}(t,s;\theta) = c\lambda^{(t+s)/2}\rho^{|t-s|}, \ \theta = [c \ \ \lambda]^T \quad (6)$$

$$k^{\mathrm{TC}}(t,s;\theta) = c\min(\lambda^t,\lambda^s), \ \theta = [c \ \ \lambda \ \ \rho]^T \quad (7)$$

$$c \geq 0, 0 \leq \lambda < 1, \begin{cases} |\rho| \leq 1, & \mathrm{DT} \\ 0 \leq \rho \leq 1, & \mathrm{CT} \end{cases}$$

where the TC kernel is a special case of the DC kernel with $\rho = \sqrt{\lambda}$ and is also called the first order stable spline kernel.

**Remark 2.1** *It is worth to note that for $\rho < 0$ and $|t-s| \notin \mathbb{N}$, $\rho^{|t-s|}$ is complex and thus $k^{\mathrm{DC}}(t,s;\theta)$ with $\rho < 0$ is not well defined for the CT case.*

### 2.4 Optimal Kernel and Stable Kernel

For the KRM, the optimal kernel in the sense of minimizing the mean square error (MSE) exists [3,4] and motivates a general guideline to design a kernel. To state this result, we assume that the data has been generated by (1) for a *true* impulse response $g^0(t)$ and we let $\bar{g}^0$ and $\hat{\bar{g}}$ to represent any finite dimensional vector obtained by sampling $g^0(t)$ and its estimate $\hat{g}(t)$ at the same arbitrary time instants. Then the following result holds.

**Lemma 2.1 (Optimal kernel, [3,4])** *Letting $\gamma = \sigma^2$. Then for the KRM, the MSE matrix*

$$MSE(k) \triangleq \mathbb{E}\left[(\hat{\bar{g}} - \bar{g}^0)(\hat{\bar{g}} - \bar{g}^0)^T\right], \quad (8)$$

*where $\mathbb{E}$ is the mathematical expectation, is minimized by the kernel*

$$k^{\mathrm{opt}}(t,s) = g^0(t)g^0(s), \quad t,s \in X \quad (9)$$

*in the sense that $MSE(k) - MSE(k^{\mathrm{opt}})$ is positive semidefinite for any positive semidefinite kernel $k$.*

The optimal kernel $k^{\mathrm{opt}}$ cannot be applied in practice as it depends on the true impulse response $g^0(t)$. However, it motivates a general *guideline* to design a kernel: *let the kernel mimic the behavior of the optimal kernel, and moreover, the prior knowledge of the true impulse response should be used in the design of the kernel.*

For instance, if the LTI system is known to be stable, i.e., $g \in \ell^1$ or $\mathcal{L}^1$, then the designed kernel $k$ should reflect this, and a necessary condition to satisfy is that $k(t,t)$ should tend to 0 as $t$ goes to infinity. This observation basically rules out the possibility to model the impulse response of stable LTI systems with stationary kernels [2] and a rigorous proof for this has been given

[2] Recall that a kernel $k(t,s)$ with $t,s \in X$ is said to be stationary if $k(t,s)$ is a function of $t-s$ for any $t,s \in X$.

in [13, Lemma 8]. More generally, the designed kernel should guarantee that its associated RKHS $\mathcal{H}_k$ is a subspace of $\ell^1$ or $\mathcal{L}^1$ and a kernel that has such property is said to be stable [13,4]. Sufficient and necessary conditions for a kernel to be stable exist and are given below.

**Lemma 2.2 ([24,13])** *A positive semidefinite kernel $k$ is stable if and only if*

$$\begin{aligned} DT: \ & \sum_{s=1}^{\infty} \left|\sum_{t=1}^{\infty} k(t,s)l(t)\right| < \infty, \ \forall \ l(t) \in \ell^{\infty} \\ CT: \ & \int_0^{\infty} \left|\int_0^{\infty} k(t,s)l(t)dt\right| ds < \infty, \ \forall \ l(t) \in \mathcal{L}^{\infty} \end{aligned} \quad (10)$$

*where $\mathcal{L}^{\infty}$ and $\ell^{\infty}$ denote the space of bounded functions on $\mathbb{R}_{\geq 0}$ and bounded sequences, respectively.*

**Corollary 2.1 ([13,4])** *A positive semidefinite kernel $k$ is stable if*

$$\begin{aligned} DT: \ & \sum_{s=1}^{\infty} \left|\sum_{t=1}^{\infty} k(t,s)\right| < \infty, \\ CT: \ & \int_0^{\infty} \left|\int_0^{\infty} k(t,s)dt\right| ds < \infty. \end{aligned} \quad (11)$$

### 2.5 Kernel Design

Hereafter, we focus on the problem of kernel design. The problem is stated as follows: for the given prior knowledge of the impulse response to be identified, our goal is to design a kernel such that the prior knowledge is embedded in the kernel. The answer should depend on the type of the prior knowledge. We will consider two types of prior knowledge and discuss how to design kernels accordingly from two perspectives: a machine learning perspective and a system theory perspective.

It should be noted that the LTI system (1) is assumed to be stable, meaning that $g \in \ell^1$ or $\mathcal{L}^1$ becomes our first prior knowledge. In this regard, Theorem 2.2 becomes a rule that the designed kernel should obey.

**Remark 2.2** *Other than stability, the cases where the impulse response or the corresponding LTI system is known to have relative degree, to be monotonic, smooth and to have delay, have been discussed in [13] and in particular, sufficient and/or necessary conditions for a kernel to have such properties are given. For instance, if a CT kernel $k$ is $m$-times continuously differentiable, then every function $h \in \mathcal{H}_k$ is $m$-times continuously differentiable [25, Corollary 4.36, p. 131], [13, Lemma 5].*

## 3 A Machine Learning Perspective

We first show that both the SS kernel and the DC kernel belong to the class of the so-called amplitude modulated locally stationary (AMLS) kernel. Accordingly, we propose to treat the impulse response as a function, whose

amplitude varies with a certain rate and decays exponentially to zero in the end, and moreover, design AMLS kernels to model the impulse response.

### 3.1 Amplitude Modulated Locally Stationary Kernels

Recall that a kernel $k(t, s)$ is said to be a locally stationary (LS) kernel in the sense of Silverman [26] if

$$k(t, s) = k^d \left( \frac{t + s}{2} \right) k^c(t - s), \quad t, s \in X \qquad (12)$$

where $k^d \geq 0$ and $k^c$ is a stationary kernel [3]. Motivated by the LS kernel, we introduce a kernel suitable for modeling impulse response, which is called the amplitude modulated locally stationary (AMLS) kernel.

**Definition 3.1** *A kernel $k(t, s)$ is said to be an amplitude modulated locally stationary kernel if*

$$k(t, s) = k^d(t, s) k^c(t - s), \ k^d(t, s) = b(t)b(s), \quad (13)$$
$$k^c(0) = 1, \quad t, s \in X,$$

*where $b(t) > 0$ is bounded and $k^c$ is a stationary kernel.*

It can be proved that the AMLS kernels have the following properties (proof given in the Appendix).

**Proposition 3.1** *Consider the AMLS kernel (13). Then the following results hold:*

*(a) $k^d(t, s)$ is a rank-1 kernel[4] and moreover satisfy*

$$k^d(t, s) = \sqrt{k^d(t, t)k^d(s, s)}. \qquad (14)$$

*(b) Let $g(t)$ be a stochastic process with zero mean and the AMLS kernel as the covariance function. Then $k^c(t - s)$ is the correlation coefficient between $g(t)$ and $g(s)$.*

*(c) Let $h(t)$ be a stochastic process with zero mean and $k^c(t - s)$ as the covariance function. Then the stochastic process $g(t) \triangleq b(t)h(t)$ has zero mean and the AMLS kernel as its covariance function.*

**Remark 3.1** *Since $k^c(t-s)$ is the correlation coefficient and $k^c(0) = 1$, then $k^c(t - s)$ is bounded and moreover, $|k^c(t - s)| \leq 1$ for any $t, s \in X$.*

For a function estimation problem, if a zero mean Gaussian processes with the AMLS kernel as the covariance function is chosen to model the function, then Proposition 3.1 has the following implications:

---

[3] Clearly, if $k^d$ is a positive constant, then the LS kernel reduces to a stationary kernel.

[4] A kernel $k(t, s)$ with $t, s \in X$ is said to be a rank-1 kernel if for any $t_i, s_i \in X$, $i = 1, \dots, n$ and for any $n \in \mathbb{N}$, the kernel matrix $K$, defined by $K_{i,j} = k(t_i, t_j)$, is a rank-1 matrix.

- The stationary kernel $k^c(t - s)$ in (13) accounts for the varying rate of the function.

  Note that part (b) shows that $k^c(t - s)$ is the correlation coefficient for $g(t)$, which implies that

  $$\mathbb{E}(g(t) - g(s))^2$$
  $$= (b(t))^2 + (b(s))^2 - 2b(t)b(s)k^c(t - s).$$

  The above equation shows that for any fixed $t, s \in X$, as $k^c(t-s)$ varies from 1 to $-1$, $(g(t) - g(s))^2$ tends to become larger, that is, $g(s)$ tends to vary more quickly away from $g(t)$. This observation also holds for $h(t)$ as $k^c(t - s)$ is also the correlation coefficient for $h(t)$.

- The rank-1 kernel $k^d(t, s)$ in (13) account for the change of the amplitude of the function.

  Note that part (c) shows that the amplitude of $g(t)$ is modulated by the factor $b(t)$. If $h(t)$ does not converge to 0 as $t$ goes to $\infty$, then a suitable $b(t)$ can be chosen such that $g(t)$ does.

Interestingly, both the SS and DC kernels can be put in the form of (13). Setting $\lambda = \exp(-\beta/2)$ and using the equality $\max\{t, s\} = (t + s + |t - s|)/2$ yields that the SS kernel (5) is rewritten as

$$k^{\mathrm{SS}}(t, s; c, \lambda) = c\lambda^{3(t+s)} \left( \frac{1}{2}\lambda^{|t-s|} - \frac{1}{6}\lambda^{3|t-s|} \right). \quad (15)$$

Then we identify, for the SS kernel (15),

$$k^d(t, s) = \frac{c}{3}\lambda^{3(t+s)}, k^c(t - s) = \frac{3}{2}\lambda^{|t-s|} - \frac{1}{2}\lambda^{3|t-s|},$$
$$(16)$$

and for the DC kernel (6),

$$k^d(t, s) = c\lambda^{\frac{1}{2}(t+s)}, k^c(t - s) = \rho^{|t-s|}. \qquad (17)$$

Moreover, one can prove the following result (proof given in the Appendix).

**Proposition 3.2** *The SS kernel (15) and the DC kernel (6) are AMLS kernels.*

Now we demonstrate for the SS and DC kernels the role of the rank-1 kernel $k^d$ and the stationary kernel $k^c$.

**Example 3.1** *As seen from the left panel of Fig. 1, as $\lambda$ changes from $0.9^{\frac{1}{2}}$ to $0.5^{\frac{1}{2}}$, the realizations of the DT zero mean Gaussian process with $k^c(t - s; [c, \lambda]^T)$ in (16) as the covariance function varies more quickly, because $k^c(1; [1, 0.9^{\frac{1}{2}}]^T) > k^c(1; [1, 0.5^{\frac{1}{2}}]^T)$. Similarly, as $\rho$ changes from 0.99, to 0 and to $-0.99$, the realizations of the DT zero mean Gaussian process with $k(t - s; [1, 0.9, \rho]^T)$ in (17) varies more quickly, and especially for the case with $\rho = -0.99$, the realization tends to change its sign at the next time instant. The above observations carry over to the right panel of Fig. 1. Finally,*
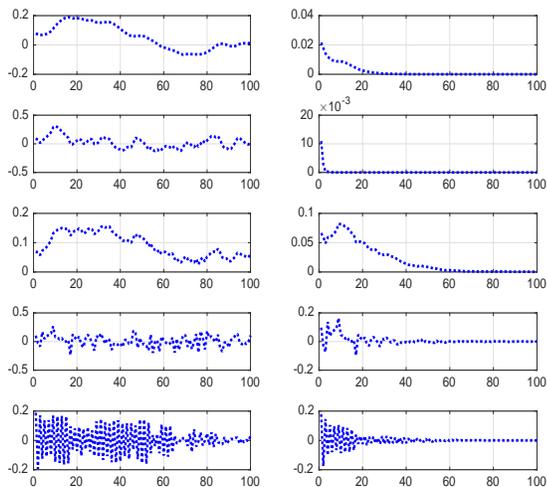
Fig. 1. Realizations of the DT zero mean Gaussian processes with the SS kernel (15) and the DC kernel (6) as covariance function. For each row, the left panel shows the realizations of the DT zero mean Gaussian process with the IS kernel $k^c(t-s)$ as covariance function and the right panel shows the realization of the DT zero mean Gaussian process with the SS or DC kernel as covariance function. The corresponding realizations for the SS kernel (15) are shown in the top two rows, which correspond to $c = 1$ and $\lambda = 0.9^{\frac{1}{2}}, 0.5^{\frac{1}{2}}$, respectively. The corresponding realizations for the DC kernel (6) are shown in the bottom three rows, which correspond to $c = 1$, $\lambda = 0.9$ and $\rho = 0.99, 0, -0.99$, respectively.

*the realizations on the left panel of Fig. 1 do not go to zero for large t but the realizations on the right panel do, because of $k^d(t, s)$.*

**Remark 3.2** *Recall from Remark 2.1 that the DC kernel (6) with $\rho < 0$ is only defined for the DT case but not for the CT case. Now we see that the DC kernel is not good to model quickly varying impulse responses for the CT case, which is also true for the SS kernel (15).*

## 3.2 Construct AMLS Kernels for Regularized Impulse Response Estimation

Proposition 3.2 and its implications on the role of the rank-1 kernel $k^d$ and the stationary kernel $k^c$ suggest a machine learning method to design kernels for regularized impulse response estimation. If the impulse response is treated as a function and the prior knowledge is about its decay and varying rate, then AMLS kernels can be designed by choosing suitable rank-1 kernel $k^d$ and stationary kernel $k^c$ to account for the decay and varying rate of the impulse response, respectively.

### 3.2.1 Stationary Kernel $k^c(t-s)$

An important class of stationary kernels is the isotropic stationary (IS) kernel. Recall that a stationary kernel $k^c(t-s)$ is said to be an IS kernel if it depends on $|t-s|$. IS kernels have been studied extensively in the literature in statistics and machine learning, see e.g. [27, Section 4.2.1]. There are many choices of IS kernels that could be used instead of the ones in (16) and (17).

Most of the IS kernels introduced in [27, pages 83 -88] decay monotonically w.r.t. $|t-s|$ and are always positive. For example, the squared exponential (SE) kernel and the Matèrn class of kernels:

$$k^c(r) = e^{-\beta r^2}, \ \beta > 0, \ \text{``SE''},$$
$$k^c(r) = \frac{2^{1-\nu}}{\Gamma(\nu)} \left(\beta\sqrt{2\nu}r\right)^\nu K^\nu \left(\beta\sqrt{2\nu}r\right),$$
$$\beta > 0, \nu > 0\text{``Matèrn''},$$

where $r = |t-s|$, $\Gamma(\cdot)$ is the Gamma function and $K^\nu$ is a modified Bessel function with order $\nu$.

There are IS kernels that do *not* decay monotonically, can take negative values, and can have the form of damped oscillation w.r.t. $|t-s|$. For example, the kernel in [27, p. 89], [28]:

$$k^c(r) = c(\alpha r)^{-\nu} J_\nu(\alpha r), \quad \alpha > 0, \nu \geq -1/2, \quad (18a)$$

where $r = |t-s|$, $c$ is a scalar such that $k^c(0) = 1$, and $J_\nu(\cdot)$ is the Bessel function of the first kind with order $\nu$. The kernel (18a) is defined for any $r \geq 0$ and can thus be used for both CT and DT impulse response estimation. In particular, for $\nu = -1/2$ and $\nu = 1/2$, the kernel (18a) takes the following form

$$k^c(r) = \cos(\alpha r), \ \alpha > 0, \quad (18b)$$
$$k^c(r) = \frac{\sin(\alpha r)}{\alpha r}, \ \alpha > 0. \quad (18c)$$

### 3.2.2 Rank-1 Kernel $k^d(t, s)$

The design of $k^d(t, s)$ is equivalent to that of the strictly positive function $b(t)$ in (13). The bottom line is that $b(t)$ should ensure that the designed AMLS kernel (13) is stable, i.e, $\mathcal{H}_k$ is a subspace of $\ell^1$ or $\mathcal{L}^1$. Our following result gives a characterization of the stability of the AMLS kernel (proof given in the Appendix).

**Proposition 3.3** *Consider the AMLS kernel (13). Then the following results hold.*

(a) *If $b(t) \in \ell^1$ for the DT case, and $b(t) \in \mathcal{L}^1$ for the CT case, then the AMLS kernel (13) is stable.*

*(b) Assume that there exists a sequence of positive numbers $\lambda_i$ and linearly independent functions $\phi_i(t)$ defined on $X$ such that*

$$k^c(t - s) = \sum_{i=1}^{\infty} \lambda_i \phi_i(t)\phi_i(s),$$
$$t, s \in X, \lambda_i > 0, i = 1, \cdots, \infty, \qquad (19)$$

*where the convergence is absolute and uniform on $Y_1 \times Y_2$ with $Y_1, Y_2$ being any compact subsets of $X$. If the AMLS kernel (13) is stable, then there exists no $\epsilon > 0$ such that $b(t) \geq \epsilon$ for all $t \in X$.*

Since the $b(t)$ in (16) and (17) are exponential decay functions and clearly satisfy $b(t) \in \ell^1$ or $\mathcal{L}^1$, the SS and DC kernels are stable by Proposition 3.3.

**Remark 3.3** *It is reasonable to check wether the condition (19) is too strong to be satisfied. In fact, the series expansion (19) can be obtained by Mercer's Theorem [29,30,31]. For instance, a sufficient condition for (19) to hold is given in [31]. To state this result, we define the kernel section of $k^c$ at a fixed $s \in X$ as $k_s^c \triangleq k^c(t - s)$ and we let $\mathcal{L}_2(X, \mu)$ denote the space of functions $f : X \to \mathbb{R}$ such that $\int_X |f(t)|^2 d\mu(t) < \infty$, where $\mu$ is a nondegenerate Boreal measure on $X$. Then (19) holds if $k^c(t - s)$ is continuous, $k_s^c \in \mathcal{L}_2(X, \mu)$ and*

$$\int_X \int_X (k^c(t - s))^2 d\mu(t) d\mu(s) < \infty. \qquad (20)$$

*It is easy to check that many stationary kernels satisfy the above sufficient condition, e.g., the SE kernel and the $k^c$ in (16) and (17).*

**Remark 3.4** *It follows from (A.6) that if there exists an $\epsilon > 0$ such that $|h(t)| \geq \epsilon$ for all $t \geq 0$, then $b(t) \in \mathcal{L}^1$, implying that $b(t) \in \mathcal{L}^1$ is also necessary for the stability of the AMLS kernel (13). The above observations show that under the assumption that the AMLS kernel (13) is stable, the result we can draw on the properties of $b(t)$ is determined by the property of $\mathcal{H}_{k^c}$.*

**Remark 3.5** *The reason why we force $b(t) > 0$ for $t \in X$ is because for a kernel in the form of $k(t, s) = k^d(t, s)k^c(t - s)$, we expect the two kernels $k^d$ and $k^c$ have somewhat independent role. As shown above, this idea eases the kernel design and the corresponding analysis. If this idea is not taken, then we can design more general $b(t)$ and even more general $k^d$. For instance, we could allow $b(t)$ to be arbitrary real-valued function and we could also allow $k^d$ to be the more general exponentially convex (EC) kernel[5] and design*

---

[5] For the LS kernel (12), if $k^d$ is also a kernel, then $k^d$ is called an exponentially convex (EC) kernel [32] and in this case $k$ is called an exponentially convex locally stationary (ECLS) kernel [26].
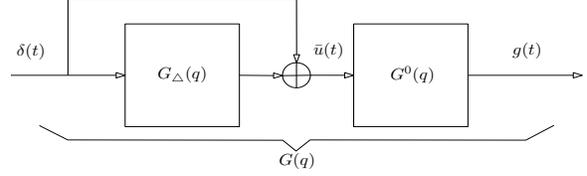
Fig. 2. The block diagram for the multiplicative uncertainty paradigm in robust control. The single-input-single-output system $G(q) = G^0(q)(1 + G_\triangle(q))$ consists of two parts: the nominal part $G^0(q)$ and the uncertainty part $G_\triangle(q)$, where $q$ is the forward shift operator and the differential operator for the DT and CT case, respectively. The real-valued signals $\delta(t), \bar{u}(t)$ and $g(t)$ are the impulsive input, the input of $G^0(q)$ and the output of $G(q)$, respectively.

*accordingly the ECLS kernel [15] instead of the AMLS kernel. The price to pay is that the role of $k^d$ and $k^c$ would become obscure and in particular, $k^d$ would not describe the decay rate and $k^c$ would not be the correlation coefficient of the underlying Gaussian process.*

## 4 A System Theory Perspective

Instead of simply treating the impulse response as a function, we now associate the impulse response with an LTI system which is stable and may be overdamped, underdamped, have multiple distinct time constants and resonant frequencies, and etc., and our goal is to design kernels to embed this kind of prior knowledge.

### 4.1 Sketch of The idea

Suppose that the prior knowledge of an LTI system is embedded in a stochastic process $g(t)$ that is used to model the corresponding impulse response. Then following the guideline to design kernels in Section 2.4 that is to mimic the behavior of the optimal kernel (9), we should design the kernel as

$$k(t, s) = \mathbb{C}ov(g(t), g(s)), \qquad (21)$$

where $\mathbb{C}ov(g(t), g(s))$ is the covariance between $g(t)$ and $g(s)$. Now the problem of kernel design becomes "how to embed the prior knowledge of an LTI system in a stochastic process $g(t)$ that is used to model the corresponding impulse response"?

A natural way to tackle the above question is by using simulation. To this goal, it is useful to employ the multiplicative uncertainty configuration in robust control, see e.g., [20], as shown in Fig. 2. Here, the nominal model $G^0(q)$ is used to embed the prior knowledge on the LTI system to be identified, the uncertainty $G_\triangle(q)$ is assumed to be stable and finally, the multiplicative uncertainty configuration is used to take into account both the nominal model and the uncertainty, and is simulated with an impulsive input to get a zero-mean Gaussian
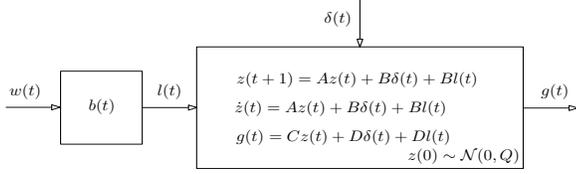
Fig. 3. The block diagram for the SI kernel (24).

process to model the impulse response $g(t)$. This idea leads to the system theory method to design kernels and we call such kernels the simulation induced (SI) kernels.

**Remark 4.1** *It is worth to note that the problem of marrying system identification and robust control is not new and has been studied over two decades in the system identification community, see e.g., [33,34]. Both the additive uncertainty configuration and the multiplicative uncertainty configuration have been used and many results have been obtained, see e.g., [35,33,34,36,37,38]. In particular, [35,33,36,37,38] used the nominal model $G_0(q)$ and the uncertainty $G_\triangle(q)$ to represent a model estimate and the corresponding model error, respectively. Here the nominal model $G_0(q)$ and the uncertainty $G_\triangle(q)$ are used to embed the prior knowledge of the LTI system and describe the corresponding uncertainty, respectively.*

### 4.2 Simulation Induced Kernel

More specifically, the prior knowledge is embedded in $G^0(q)$ or equivalently, in its state-space model

$$G^0(q): \quad \begin{aligned} z(t+1) &= Az(t) + B\bar{u}(t) && \text{(22a)} \\ \dot{z}(t) &= Az(t) + B\bar{u}(t) && \text{(22b)} \\ z(0) &\sim \mathcal{N}(0, Q) && \text{(22c)} \\ g(t) &= Cz(t) + D\bar{u}(t) && \text{(22d)} \end{aligned}$$

where $A, B, C, D,$ and $Q$ have compatible dimensions, $\bar{u}(t)$ and $g(t)$ are the input and output of $G^0(q)$, respectively. In what follows, we will use the quintuple $(A, B, C, D, Q)$ to represent the state-space model (22).

Since the uncertainty $G_\triangle(q)$ is stable, the simplest way to model the impulse response of $G_\triangle(q)$ with a stochastic process $l(t)$, is to have

$$G_\triangle(q): \quad l(t) = b(t)w(t), \quad t \in X, \qquad \text{(23)}$$

where $w(t)$ is a white Gaussian noise with zero mean and unit variance (independent of $z(0)$), and $b(t) > 0$ and $b(t) \in \ell^1$ or $\mathcal{L}^1$ for the DT and CT case, respectively. Clearly, $l(t)$ is a zero mean Gaussian process with the AMLS kernel $b(t)b(s)\delta(t-s)$ as the covariance function, where $\delta(\cdot)$ is the Dirac Delta function and the AMLS kernel $b(t)b(s)\delta(t-s)$ is stable by Proposition 3.3.

Finally, we simulate the system $G(q)$ in Fig. 2 with the impulsive input $\delta(t)$ and noting (22), (23) and $\bar{u}(t) =$

$\delta(t) + l(t)$, we obtain the simulation induced (SI) kernel:

$$\begin{aligned} z(t+1) &= Az(t) + B\delta(t) + Bb(t)w(t), && \text{(24a)} \\ \dot{z}(t) &= Az(t) + B\delta(t) + Bb(t)w(t), && \text{(24b)} \\ g(t) &= Cz(t) + D\delta(t) + Db(t)w(t), && \text{(24c)} \\ z(0) &\sim \mathcal{N}(0, Q) && \text{(24d)} \\ k^{\text{SI}}(t,s) &= \mathbb{C}ov(g(t), g(s)), && \text{(24e)} \end{aligned}$$

whose block diagram is shown in Fig. 3. From linear system theory, e.g., [39], the formal expression of the SI kernel (24) is available. For the DT case, it is

$$\begin{aligned} k^{\text{SI}}(t,s) &= CA^t Q (A^s)^T C^T + D^2 b(t) b(s) \delta(t-s) \\ &+ Db(t) \sum_{k=0}^{s-1} \delta(t-k) b(k) B^T (A^{s-1-k})^T C^T \\ &+ Db(s) \sum_{k=0}^{t-1} \delta(s-k) b(k) B^T (A^{t-1-k})^T C^T \\ &+ \sum_{k=0}^{\min\{t,s\}-1} b(k)^2 CA^{t-1-k} BB^T (A^{s-1-k})^T C^T \end{aligned} \tag{25a}$$

For the CT case with $D = 0$, it is

$$\begin{aligned} k^{\text{SI}}(t,s) &= Ce^{At} Q (e^{As})^T C^T + \\ C &\int_0^{\min\{t,s\}} b^2(\tau) e^{A(t-\tau)} BB^T (e^{A(s-\tau)})^T d\tau C^T \end{aligned} \tag{25b}$$

**Remark 4.2** *It should be noted that both DT and CT cases are considered in (22) and (24). For the CT case, (24b) could be more rigorously written as an Itô stochastic differential equation (SDE) as in [14, eq. (9)]. However, we decide to take the same point of view as in [40] to use (24b) instead in order to save the space for the introduction of SDE stuff, which is only used to derive (25b).*

**Remark 4.3** *The SI kernels (25) may or may not have closed form expressions. The related computational difficulty and cost depend on whether or not the SI kernels (25) have closed form expressions. If they do, then the computation would be easier and similar to the SS and DC kernels. If they do not, the computation of the hyperparameter estimate and the regularized impulse response estimate would become more demanding. In this regard, the particle filtering based technique for nonlinear state-space model identification in [41] could be adopted. However, the technique in [41] cannot be applied trivially as the measurement output (1) for the state space model (24) does not depend on the current state $z(t)$ solely but also the past state due to the presence of convolution. More details will be reported in an independent paper.*

Note that $g(t)$ in (24) is a Gaussian process with the impulse response of the LTI system $(A, B, C, D, Q)$ as its mean and the SI kernel $k^{\text{SI}}(t,s)$ as its covariance

function. If $\delta(t)$ is set to 0 in (24), $g(t)$ in (24) is a zero mean Gaussian process with the SI kernel $k^{\mathrm{SI}}(t,s)$ as its covariance function. In what follows, when considering kernel design, we will set $\delta(t)$ to zero for convenience.

Interestingly, the AMLS kernel (13) is closely related with the SI kernel (24), and in fact many AMLS kernels such as the SS and DC kernels can be put in the form of (24) and are thus SI kernels. To state this result, recall that the power spectral density of $k^c(t-s)$, denoted by $\Psi(\omega)$, is defined as

$$\Psi(\omega) = \begin{cases} \sum_{\tau=-\infty}^{+\infty} k^c(\tau)e^{-i\omega\tau}, & \mathrm{DT}, \\ \int_{-\infty}^{+\infty} k^c(\tau)e^{-i\omega\tau}d\tau, & \mathrm{CT}. \end{cases} \quad (26)$$

**Proposition 4.1** *Consider the AMLS kernel (13). Assume that $k^d(t,s) = c\lambda^{t+s}$ with $c \geq 0$ and $0 \leq \lambda < 1$ and $\Psi(\omega)$ is a proper rational function of $e^{i\omega}$ or $\cos(\omega)$ for the DT case, and $\omega$ for the CT case, respectively. Then the AMLS kernel (13) can be put in the form of (24) and thus is SI kernel.*

**Proposition 4.2** *Consider the SS kernel (15) and the DC kernel (6). Then the following results hold:*

- *For the DT case, let*

$$\bar{a} = \sqrt{1 + \lambda^2 - \sqrt{1 + \lambda^2 + \lambda^4}},$$
$$\bar{b} = \sqrt{1 + \lambda^2 + \sqrt{1 + \lambda^2 + \lambda^4}}.$$

*Then the SS kernel is in the form of (24) with*

$$A = \lambda^3 \begin{bmatrix} \lambda & 0 \\ 0 & \lambda^3 \end{bmatrix}, B = \begin{bmatrix} \lambda^3 \\ \lambda^3 \end{bmatrix}, Q = \frac{c}{3} \begin{bmatrix} \frac{1}{1-\lambda^2} & \frac{1}{1-\lambda^4} \\ \frac{1}{1-\lambda^4} & \frac{1}{1-\lambda^6} \end{bmatrix}$$
$$C = \sqrt{\frac{(1-\lambda^2)^3}{2}} \begin{bmatrix} \frac{\bar{a}+\bar{b}\lambda}{1-\lambda^2} & \frac{-\lambda^2(\bar{a}+\bar{b}\lambda^3)}{1-\lambda^2} \end{bmatrix}$$
$$D = \bar{b}\sqrt{\frac{(1-\lambda^2)^3}{2}}, b(t) = \left(\frac{c}{3}\right)^{\frac{1}{2}}\lambda^{3t}. \quad (27)$$

*The DC kernel is in the form of (24) with*

$$A = \lambda^{\frac{1}{2}}\rho, B = \lambda^{\frac{1}{2}}, C = \rho(1-\rho^2)^{\frac{1}{2}}, Q = \frac{c}{1-\rho^2},$$
$$D = (1-\rho^2)^{\frac{1}{2}}, b(t) = c^{\frac{1}{2}}\lambda^{\frac{t}{2}}. \quad (28)$$

- *For the CT case, let $\lambda = \exp(-\beta/2)$ and $\rho = \exp(-\alpha)$. Then the SS kernel is in the form of (24) with*

$$A = \begin{bmatrix} -\frac{3}{2}\beta & 1 \\ -\frac{3}{4}\beta^2 & -\frac{7}{2}\beta \end{bmatrix}, B = \begin{bmatrix} 0 \\ 1 \end{bmatrix}, Q = \frac{c}{3}\begin{bmatrix} \frac{1}{3\beta^3} & 0 \\ 0 & \frac{1}{4\beta} \end{bmatrix}$$
$$C = \begin{bmatrix} 3^{\frac{1}{2}}\beta^{\frac{3}{2}} & 0 \end{bmatrix}, D = 0, b(t) = \left(\frac{c}{3}\right)^{\frac{1}{2}}e^{-\frac{3}{2}\beta t}. \quad (29)$$

*The DC kernel is in the form of (24) with*

$$A = -\alpha - \frac{1}{4}\beta, B = 1, C = (2\alpha)^{\frac{1}{2}}, Q = \frac{c}{2\alpha},$$
$$D = 0, b(t) = c^{\frac{1}{2}}e^{-\frac{\beta}{4}t}. \quad (30)$$

Proposition 4.2 enhances our understanding about the SS and DC kernels and in particular, the prior knowledge embedded in the two kernels from a system theory perspective is as follows:

1) For the SS kernel, the corresponding nominal model $G^0(q)$ in Fig. 2 is a second order system. For the DC kernel, the corresponding nominal model $G^0(q)$ is a first order system. For both the SS and DC kernels and for the CT case, $G^0(q)$ has negative real poles corresponding to overdamped impulse responses. For the DT case, the two poles for the SS kernel are real positive and correspond to impulse response without oscillation, and the pole for the DC kernel can be positive or negative (depending on $\rho$) and corresponds to impulse response without or with oscillation. The pole for the TC kernel is positive (equal to $\lambda$) and corresponds to impulse response without oscillation.
2) For both the SS and DC kernels, the decay rate of the impulse response of the uncertainty $G_\triangle(q)$ in Fig. 2 are all described by exponential decay functions.

**Remark 4.4** *It is worth to note that to study the link between a kernel and its state-space model realization has a long history, see e.g., [42,40]. This link is important because it opens the body of a kernel from a system theory perspective and accordingly helps to understand the underlying behavior of a kernel.*

### 4.3 Stability of SI kernels

For the SI kernel (25), we provide the following sufficient condition to guarantee its stability.

**Proposition 4.3** *Consider the SI kernel (25).*

- *For the DT case, assume that $A$ has distinct eigenvalues and moreover, $A$ is stable, i.e., $A$ has all eigenvalues strictly inside the unit circle. Assume further that*

$$b(t) \leq \bar{c}|\bar{\lambda}|^{\frac{t}{2}} \quad (31a)$$

*for some $\bar{c} > 0$, where $\bar{\lambda}$ is the eigenvalue of $A$ with the largest modulus. Then the SI kernel (25a) is stable.*
- *For the CT case, assume that $A$ has distinct eigenvalues and moreover $A$ is stable, i.e., $A$ has all eigenvalues with strictly negative real parts. Assume further that*

$$b(t) \leq \bar{c}e^{-\frac{1}{2}|\mathrm{Re}(\bar{\lambda})|t} \quad (31b)$$

9

for some $\bar{c} > 0$, where $\bar{\lambda}$ is the eigenvalue of $A$ with the largest real part and $\mathrm{Re}(\bar{\lambda})$ is the real part of $\bar{\lambda}$. Then the SI kernel (25b) is stable.

**Remark 4.5** *The conditions (31) are sufficient but not necessary. For instance, consider the DT case. The SS kernel (27) satisfies (31a) but the DC kernel (28) does not. It is possible to derive a slower upper bound for $b(t)$. Another intension to have Proposition 4.3 is to give some guidelines about the relations between $b(t)$, $A$ and the stability of SI kernel: in order to guarantee stability of the SI kernel, the impulse response of the uncertainty $G_\triangle(q)$ (23), i.e., $b(t)$ should decay a bit faster than the slowest mode of the nominal model $G^0(q)$.*

### 4.4 Maximum Entropy Property of SI Kernels

Since the prior knowledge is never complete, it is worth to note Jaynes's maximum entropy (MaxEnt) rationale [43] to derive from incomplete prior knowledge the optimal statistical prior distribution. Jaynes's idea is to formulate a MaxEnt problem with respect to the prior, and then solve the problem to obtain the optimal prior. The constraints of the problem describes the prior knowledge (in the MaxEnt sense) of the underlying stochastic process (the system) to be identified. Interestingly, the SI kernel lends itself easily to a MaxEnt interpretation, leading to a new facet to understand the underlying behavior of the SI kernel.

Recall that the differential entropy $H(X)$ of a real-valued continuous random variable $X$ is defined as $H(X) = -\int_S p(x) \log p(x) dx$, where $p(x)$ is the probability density function of $X$ and $S$ is the support set of $X$. Then we prove the next MaxEnt interpretation of SI kernels[6].

**Proposition 4.4** *Consider the SI kernel (24) with (25a).*

- *For the case $D \neq 0$, define*

$$f(t) = \frac{\bar{g}(t) - CA^t\bar{z}(0) - \sum_{k=0}^{t-1} CA^{t-1-k}Bb(k)f(k)}{Db(t)}$$
$$\bar{g}(t) \in \mathbb{R}, \bar{z}(0) \in \mathbb{R}^n, t = 0, 1, \cdots, s \qquad (32)$$

*Then for any $s \in \mathbb{N}$, the Gaussian process $g(t)$ in (24) is the solution to the MaxEnt problem*

$$\underset{\bar{z}(0),\bar{g}(t)}{\text{maximize}} \quad H(\bar{z}(0), \bar{g}(0), \bar{g}(1), \cdots, \bar{g}(s)) \qquad (33a)$$
$$\text{subject to}$$
$$\mathbb{E}(\bar{z}(0)) = 0, \mathbb{E}(\bar{g}(t)) = 0,$$
$$\mathbb{C}ov(\bar{z}(0)) = Q, \mathbb{V}(f(t)) = 1, t = 0, \cdots, s \quad (33b)$$

---
[6] Let $X_1, X_2$ be two jointly distributed random variables. Then the joint differential entropy of $H([X_1, X_2]^T)$ is simply written as $H(X_1, X_2)$ below.

where $\mathbb{C}ov(\bar{z}(0))$ is the covariance matrix of $z(0)$ and $\mathbb{V}(f(t))$ is the variance of $f(t)$.

- *For the case $D = 0$, assume that $CB \neq 0$ and define*

$$f(t-1) = \frac{\bar{g}(t) - CA^t\bar{z}(0) - \sum_{k=0}^{t-2} CA^{t-1-k}Bb(k)f(k)}{CBb(t-1)}$$
$$\bar{g}(t) \in \mathbb{R}, \bar{z}(0) \in \mathbb{R}^n, t = 0, \cdots, s - 1$$
$$(34)$$

*Then for any $s \in \mathbb{N}$, the Gaussian process $g(t)$ in (24) is the solution to the MaxEnt problem*

$$\underset{\bar{z}(0),\bar{g}(t)}{\text{maximize}} \quad H(\bar{z}(0), \bar{g}(1), \cdots, \bar{g}(s)) \qquad (35a)$$
$$\text{subject to}$$
$$\mathbb{E}(\bar{z}(0)) = 0, \mathbb{E}(\bar{g}(t)) = 0,$$
$$\mathbb{C}ov(\bar{z}(0)) = Q, \mathbb{V}(f(t)) = 1, t = 0, \cdots, s - 1$$
$$(35b)$$

**Corollary 4.1** *For any $s \in \mathbb{N}$, the zero mean Gaussian process $g(t)$ with the DC kernel (6) defined on $\mathbb{N} \times \mathbb{N}$ as its covariance is the solution to the MaxEnt problem*

$$\underset{\bar{g}(t)}{\text{maximize}} \quad H(\bar{g}(0), \bar{g}(1), \cdots, \bar{g}(s))$$
$$\text{subject to} \quad \mathbb{E}(\bar{g}(t)) = 0, t = 0, \cdots, s, \mathbb{V}(\bar{g}(0)) = c,$$
$$\mathbb{V}(\bar{g}(t) - \lambda^{1/2}\rho\bar{g}(t-1)) = c(1-\rho^2)\lambda^t, t = 1, \cdots, s$$
$$(36)$$

**Remark 4.6** *When $\rho = \lambda^{1/2}$, the DC kernel (6) becomes the TC kernel (7). For the TC kernel defined on $\mathbb{N} \times \mathbb{N}$, (36) becomes*

$$\underset{\bar{g}(t)}{\text{maximize}} \quad H(\bar{g}(0), \bar{g}(1), \cdots, \bar{g}(s))$$
$$\text{subject to} \quad \mathbb{E}(\bar{g}(t)) = 0, t = 0, \cdots, s, \mathbb{V}(\bar{g}(0)) = c,$$
$$\mathbb{V}(\bar{g}(t) - \lambda\bar{g}(t-1)) = c(1-\lambda)\lambda^t, t = 1, \cdots, s$$
$$(37)$$

*which is different from the MaxEnt interpretation in [17]. It shall be noted that both interpretation are correct but derived in different ways: the MaxEnt problems are different but have the same optimal solution.*

**Remark 4.7** *By using Corollary 4.1 and the trick in [17, Theorem 2], it is possible to derive a more concise proof for [18, Proposition IV.1] which shows that, the fact that the DC kernel has tridiagonal inverse can be given a MaxEnt covariance completion interpretation.*

**Remark 4.8** *The CT case is not discussed here because according to our best knowledge the entropy is not well defined for CT stochastic processes.*

### 4.5 Markov Property of SI Kernels

As shown in [44,19,18,17], the kernel matrix of DC kernel (6) has tridiagonal inverse. Here, we further show that

the Gaussian process with the DC kernel as its covariance function is also a Markov process with order 1 and moreover, we are able to design SI kernels which correspond to more general Gaussian Markov processes and have banded [7] inverses of their kernel matrices.

First, recall from e.g., [27, Appendix B] that a Gaussian Markov process is a stochastic process that is both a Gaussian process and a Markov process. A well-known instance is the DT autoregressive process of order $p$:

$$x(t+1) = \sum_{k=0}^{p-1} a(t,k)x(t-k) + b(t+1)w(t+1) \quad (38a)$$

$$\text{or, } x(t+1) = \sum_{k=0}^{p-1} a(t,k)x(t-k) + b(t)w(t) \quad (38b)$$

where $t \in \mathbb{N}$, $x(t), a(t,k), b(t) \in \mathbb{R}$, $x(0)$ is Gaussian distributed and assumed to be independent of $w(t)$, a zero mean white Gaussian noise with unit variance. The stochastic process (38) is a Markov process with order $p$ since $x(t+p+1)$ only depends on $x(t+p), \cdots, x(t+1)$ given the history $x(s)$ with $s \leq t+p$.

**Proposition 4.5** *Consider a zero mean DT Gaussian process $g(t)$ with the DC kernel (6) as its covariance. Then $g(t)$ with $t \in \mathbb{N}$ can be put in the form of*

$$g(t+1) = \lambda^{\frac{1}{2}}\rho g(t) + (1-\rho^2)^{\frac{1}{2}} c^{\frac{1}{2}} \lambda^{\frac{t+1}{2}} w(t+1) \quad (39)$$

*and thus a Markov process with order 1, and moreover, its kernel matrix has a 1-band matrix inverse.*

**Remark 4.9** *Interestingly, (39) can also be derived from (36), i.e., from the MaxEnt property of the DC kernel.*

It is possible to construct more general SI kernels with Markov property.

**Proposition 4.6** *Consider the DT SI kernel (24) with $(A, B, C, D)$ being a realization of $G^0(q)$ which is an $n$th order DT system*

$$G^0(q) = \frac{q^{n-1}\bar{b}}{q^n + a_1 q^{n-1} + \cdots + a_n}, \quad (40)$$

*where $\bar{b}, a_1, \ldots, a_n \in \mathbb{R}$. Then the Gaussian process $g(t)$ in (24) with any $b(t) > 0$ and positive semidefinite $Q$ is also a Markov process with order $n$ and thus the SI kernel has an $n$-band matrix inverse.*

For illustration, we consider an example.

---

[7] A real symmetric matrix $A$ with dimension $n > m+1$ is called an $m-band$ matrix if $A_{i,j} = 0$ for $|i-j| > m$.
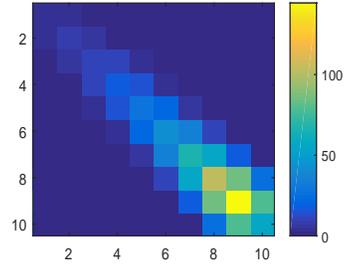


Fig. 4. Scaled image of $K^{-1}$ in Example 4.1. When generating $K$, we choose $\bar{b} = 1$, $a_1 = 0.5$, $a_2 = 0.9$, $b(t) = 0.8^t$, and $Q = I_2$. The image is drawn by using `imagesc` in MATLAB, where the colder the color the smaller the element of $K^{-1}$.

**Example 4.1** *Consider the DT SI kernel (24) with $(A, B, C, D)$ being a realization of $G^0(q)$, which is a 2nd order DT system having two stable real poles, i.e.,*

$$G^0(q) = \frac{\bar{b}q}{(q+a_1)(q+a_2)}, \quad (41)$$

*where $\bar{b} \in \mathbb{R}$ and $|a_i| < 1, i = 1, 2$. We consider the inverse of the kernel matrix $K$ defined by $K_{i,j} = k^{SI}(i,j)$ with $i, j = 1, \ldots, 10$. By Proposition 4.6, $K^{-1}$ should be a 2-band matrix, which is confirmed by Fig. 4.*

**Remark 4.10** *The Markov property of a kernel and its associated special structure (the banded inverse of the kernel matrix) can be used to develop numerically more stable and efficient implementations for this kernel based regularization method, see e.g., [18, Section 5].*

**Remark 4.11** *The CT case is not discussed here because according to [27, p. 217], regular sampling of a CT Gaussian Markov process entropy in general would not lead to a DT Gaussian Markov process. That is to say, even if a CT SI kernel with Markov property is constructed, its corresponding kernel matrix evaluated on the sampling instants may not have banded matrix inverse.*

## 5 A Case Study: Impulse Response with Damped Oscillation

In this section, we consider the estimation of impulse responses with damped oscillation and demonstrate how to design kernels from the proposed two perspectives.

### 5.1 The Machine Learning Perspective

As shown in Section 3, the machine learning perspective treats the impulse response as a function, and designs AMLS kernels with the rank-1 kernel and the stationary kernel to account for the decay and varying rate of the impulse response, respectively. Now we show that by further exploiting the rank-1 kernel or the stationary kernel, we can design AMLS kernels capable to embed
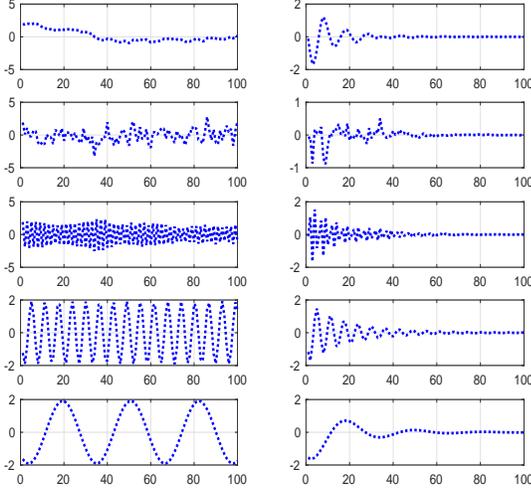
Fig. 5. Realizations of DT zero mean Gaussian process with AMLS kernels (43). For each row, the left panel shows the realizations of the DT zero mean Gaussian process with the IS kernel $k^c(t-s)$ as covariance function and the right panel shows the realization of the DT zero mean Gaussian process with the AMLS kernel as covariance function. The top three rows correspond to the AMLS kernel (43b) with $c = 1$, $\lambda = 0.9^{\frac{1}{2}}$, $\omega = 0.2\pi$ and $\rho = 0.99, 0, -0.99$, respectively. The bottom two rows correspond to the AMLS kernel (43a) with $c = 1$, $\lambda = 0.9^{\frac{1}{2}}$ and $\alpha = \pi, 0.2\pi$, respectively.

the extra prior knowledge that the impulse response has damped oscillation.

To have oscillation behavior in the AMLS kernel, we can choose to have it either in the stationary kernel or in the rank-1 kernel, i.e., $b(t)$. For example, we can choose $b(t)$ as an exponential decay function, i.e., $b(t) = c\lambda^t$ and then choose (18b) as the stationary kernel, because (18b) has oscillation behavior. Or we can choose the stationary kernel in (17) and $b(t)$ as an exponential decay function with oscillation, i.e.,

$$b(t) = c^{\frac{1}{2}}\lambda^t[\cos(\omega t) + 1 + \epsilon], \qquad (42)$$

where $c > 0$, $0 \le \lambda < 1$, $\omega \ge 0$, and $\epsilon > 0$ is a tiny number to ensure that $b(t) > 0$ for $t \in X$. The above idea leads to the following two kernels, respectively:

$$k^{\text{AMLS-2Os}}(t, s) = c\lambda^{t+s}\cos(\alpha|t-s|), \qquad (43a)$$
$$k^{\text{AMLS-2Od}}(t, s) = c\lambda^{t+s}[\cos(\omega t) + 1 + \epsilon]$$
$$\times [\cos(\omega s) + 1 + \epsilon]\rho^{|t-s|}, \quad (43b)$$

Fig. 5 illustrates that the AMLS kernels (43) are capable to describe functions with damped oscillation behavior.

## 5.2 The System Theory Perspective

As shown in Section 4, the system theory perspective associates the impulse response with an LTI system, and designs SI kernels with the nominal model to embed the prior knowledge on the LTI system. Now we show that by choosing the nominal model to be a second order LTI system with a pair of complex conjugate poles, we design a SI kernel capable to model impulse responses with damped oscillation.

More specifically, we choose the transfer function of the nominal model $G^0(q)$ in Fig. 2 to be

$$G^0(s) = \frac{1}{s^2 + 2w_0\xi s + w_0^2} \qquad (44)$$

where $\omega_0 > 0$ and $0 \le \xi < 1$. By setting $\alpha = w_0\xi$ and $\beta = w_0\sqrt{1-\xi^2}$, a state space model for (44) is described by $(A, B, C, D)$ with

$$A = \begin{bmatrix} 0 & 1 \\ -\alpha^2 - \beta^2 & -2\alpha \end{bmatrix}, B = \begin{bmatrix} 0 \\ 1 \end{bmatrix}, C = \begin{bmatrix} 1 & 0 \end{bmatrix}, D = 0.$$

Finally, setting $Q = I_2$ and $b(t) = e^{-\gamma t}$ with $\gamma > 0$ yields that the SI kernel (25b) takes the form:

$$k^{\text{SI-2Od}}(t, s) = e^{-\alpha(t+s)}[\cos(\beta t)\cos(\beta s)$$
$$+ \frac{\alpha}{\beta}\sin(\beta(t+s)) + \frac{\alpha^2+1}{\beta^2}\sin(\beta t)\sin(\beta s)]$$
$$+ \frac{e^{-\alpha(t+s)}\cos(\beta(t-s))(e^{2(\alpha-\gamma)\min\{t,s\}} - 1)}{4\beta^2(\alpha-\gamma)}$$
$$+ \frac{e^{-\alpha(t+s)}}{2\beta^2\sqrt{4\beta^2 + 4(\alpha-\gamma)^2}} \times [\cos(\phi + \beta(t+s)) -$$
$$e^{2(\alpha-\gamma)\min\{t,s\}}\cos(2\beta\min\{t,s\} - \phi - \beta(t+s))],$$
$$\phi = \arccos\left(\frac{2(\alpha-\gamma)}{4\beta^2 + 4(\alpha-\gamma)^2}\right). \qquad (45)$$

## 5.3 Numerical Simulation

To illustrate that the AMLS kernls (43) and the SI kernel (45) are capable to model LTI systems with strong oscillation, we consider the following numerical example.

### 5.3.1 Test Data-bank

The way in [45] is used to generate the test systems and data bank. In particular, we first generate 200 test systems with strong oscillation:

$$G(q) = \frac{q + 0.99}{q} \sum_{i=1}^{N_r+1} G_i(q) \qquad (46)$$

12

where

$$G_i(q) = K_i \frac{q + 0.9}{(q - p_i)(q - p_i^*)}, \ i = 1, \dots, N_r,$$

and $G_{N_r+1}(q)$ is a 4th order system randomly generated by the function `drmodel` in MATLAB with its poles inside the disk of radius 0.95. The parameters $N_r, p_i, K_i$ are randomly generated as follows: $N_r \sim \mathcal{U}[3, 8]$, $K_i \sim \mathcal{U}[2, 10]$, $p_i = \rho_i e^{j[\phi_0 + \frac{\pi}{2N_r}(i-1)]}$ with $\phi_0 \sim \mathcal{U}[0, \frac{\pi}{2}]$ and $\rho_i \sim \mathcal{U}[0.8, 0.99]$.

For each of the 200 DT systems (46), we generated the test data as follows. The function `idinput` in MATLAB is used to generate a random Gaussian input $u(t)$ with normalized band $[0, 0.95]$ and length 210. The chosen DT system is simulated with $u(t)$ to get the noise-free output $G(q)u(t)$, and $G(q)u(t)$ is then perturbed by a white Gaussian noise, whose standard deviation is equal to $\mathcal{U}[0.5, 1]$ times the standard deviation of $G(q)u(t)$.

In this way, we generate 200 test DT systems and for each system a data set with 210 data points.

### 5.3.2 Simulation Results

The following model fit is used to measure how good the estimated impulse response is:

$$fit = 100 \left( 1 - \left[ \frac{\sum_{k=1}^{100} |g_k^0 - \hat{g}_k|^2}{\sum_{k=1}^{100} |g_k^0 - \bar{g}^0|^2} \right]^{\frac{1}{2}} \right), \ \bar{g}^0 = \frac{1}{100} \sum_{k=1}^{100} g_k^0$$

where $g_k^0$ and $\hat{g}_k$ are the true impulse response and its estimate at the $k$th time instant, respectively.

The AMLS kernels (43) and the SI kernel (45) are compared with the TC kernel (7), the DC kernel (6) and the SS kernel (15) enriched with a second order parametric part proposed in [6] and denoted by SSp below.

The simulation result is summarized below. The average model fits for the tested kernels are shown below:

| | TC | DC | SSp | AMLS-2Os | AMLS-2Od | SI-2Od |
|---|---|---|---|---|---|---|
| Avg. Fit | 47.5 | 50.0 | 50.9 | 48.4 | 49.7 | 53.3 |

The distribution of the model fits are shown in Fig. 6.

### 5.3.3 Findings

For the test systems and data bank, we have the following findings. The SI kernel (45) is best among all the tested kernels in both the average accuracy and robustness and thus is best to model systems with strong oscillation. The AMLS kernel (43b) works well as it is
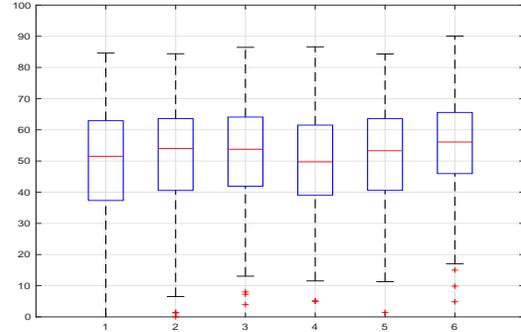


Fig. 6. Boxplot of the model fits for the TC, DC, SSp, AMLS-2Os, AMLS-2Od and SI-2Od kernels [from left to right]. The six kernels have 5, 3, 3, 3, 4, and 3 fits smaller than zero, respectively, which are not shown in the boxplot for better display.

better than the TC kernel in both the average accuracy and robustness and it is just a bit worse than the DC kernel in the average accuracy. The AMLS kernel (43a) works not very well, though it has better average accuracy than the TC kernel. One possible explanation why the SI kernel (45) is best to model impulse response with strong oscillation is that the prior knowledge embedded in the SI kernel (45) is closest to the truth. This may be true because it makes use of the prior knowledge from a system theory perspective directly by choosing the nominal model to be a second order system with a pair of complex conjugate poles.

**Remark 5.1** *To embed the prior knowledge such as multiple distinct time constants and multiple resonant frequencies, we proposed to use the multiple kernel introduced in [7] where the multiple kernel is a conic combination of some fixed kernels. The fixed kernels can be instances of the SS, TC, DC and SI kernels (24) evaluated on a grid in their hyper-parameter space. The interested readers are referred to [7] for more details.*

## 6 Conclusion

Kernel methods or regularization methods to estimate the impulse response of a linear time invariant systems is a most useful relatively recent addition to the techniques of system identification. Earlier results have demonstrated that important improvements in estimation quality can be achieved with kernels devised in a more or less ad hoc way. In this contribution, we have focused on systematic mechanisms and design concepts for how to construct kernels that are capable of adjusting its hyperparameters to capture the unknown system's properties in useful ways. The two main avenues to such thinking have been a *machine learning perspective* focusing on function properties of the impulse response and a *system theory perspective* focusing on the LTI system that produces the impulse response. This has lead

to understanding and insights into general properties of the earlier suggested methods, e.g. that the so called DC and SS kernels (derived from more ad hoc thinking) belong to the class of amplitude modulated locally stationary kernels (Proposition 3.2) and that they are simulation induced from certain LTI systems (Proposition 4.2). They are also related to maximum entropy optimal choices (Proposition 4.4) which is a valuable feature.

The take-home message of this contribution is as follows. The issue of kernel design should relate to the type of the prior knowledge and different types of prior knowledge should lead to different ways to design the corresponding kernels. Here, a machine learning perspective and a system theory perspective are introduced accordingly:

- Machine Learning perspective: If the impulse response is treated as a function and the prior knowledge is on its decay and varying rate, then we can design the amplitude modulated locally stationary (AMLS) kernel. In particular, we design a rank-1 kernel and a stationary kernel to account for the decay and varying rate of the impulse response, respectively. Moreover, by further exploiting the rank-1 kernel or the stationary kernel, it is possible to design AMLS kernels capable to embed more general prior knowledge.

- System Theory perspective: If the impulse is associated with an LTI system and the prior knowledge is that the LTI system is stable and may be overdamped, underdamped, have multiple distinct time constants and resonant frequencies and etc., then we can design the simulation induced (SI) kernel. In particular, the nominal model is used to embed the prior knowledge, the uncertainty is assumed to be stable and finally, the multiplicative uncertainty configuration is used to take into account both the nominal model and the uncertainty, and is simulated with an impulsive input to get the SI kernel or equivalently the zero-mean Gaussian process to model the impulse response.

Finally, finding a suitable kernel structure is only one leg of the kernel-based regularization method. Tuning its hyperparameters regardless of structure is the other main topic. This has been discussed in some detail e.g. in [4] and [8]. Some related new asymptotic results are recently presented in [9].

## Acknowledgements

## References

[1] L. Ljung. *System Identification - Theory for the User*. Prentice-Hall, Upper Saddle River, N.J., 2nd edition, 1999.

[2] T. Söderström and P. Stoica. *System Identification*. Prentice-Hall Int., London, 1989.

[3] T. Chen, H. Ohlsson, and L. Ljung. On the estimation of transfer functions, regularizations and Gaussian processes - Revisited. *Automatica*, 48:1525–1535, 2012.

[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*, 50(3):657–682, 2014.

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

[6] G. Pillonetto, A. Chiuso, and G. De Nicolao. Prediction error identification of linear systems: a nonparametric Gaussian regression approach. *Automatica*, 47(2):291–305, 2011.

[7] 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*, (11):2933–2945, 2014.

[8] Gianluigi Pillonetto and Alessandro Chiuso. Tuning complexity in regularized kernel-based regression and linear system identification: The robustness of the marginal likelihood estimator. *Automatica*, 51:106 – 117, 2015.

[9] B. Mu, T. Chen, and L. Ljung. Tuning of hyperparameters for fir models – an asymptotic theory. In *Proceedings of the IFAC 2017 World Congress*, page under review, Toulouse, France., 2017.

[10] A. N. Tikhonov and V. Y. Arsenin. *Solutions of Ill-posed Problems*. Winston/Wiley, Washington, D.C., 1977.

[11] J Sjöberg, T McKelvey, and L Ljung. On the use of regularization in system identification. In *Proceedings of the IFAC 2012 World Congress*, pages 381–386, Sydney, Australia., 1993.

[12] A. Chiuso. Regularization and Bayesian learning in dynamical systems: Past, present and future. *Annual Reviews in Control*, 41:24 – 38, 2016.

[13] F. Dinuzzo. Kernels for linear time invariant system identification. *SIAM Journal on Control and Optimization*, 53(5):3299–3317, 2015.

[14] Tianshi Chen and Lennart Ljung. Constructive state-space model induced kernels for regularized system identification. In *19th IFAC World Congress*, pages 1047–1052, Cape town, South Africa, 2014.

[15] T. Chen and L. Ljung. On kernel structure for regularized system identification (i): a machine learning perspective. In *Proceedings of the IFAC Symposium on System Identification*, pages 1035–1040, Beijing, China., 2015.

[16] T. Chen and L. Ljung. On kernel structure for regularized system identification (ii): a system theory perspective. In *Proceedings of the IFAC Symposium on System Identification*, pages 1041–1046, Beijing, China., 2015.

[17] Tianshi Chen, Tohid Ardeshiri, Francesca P. Carli, Alessandro Chiuso, Lennart Ljung, and Gianluigi Pillonetto. Maximum entropy properties of discrete-time first-order stable spline kernel. *Automatica*, 66:34 – 38, 2016.

[18] F. P. Carli, T. Chen, and L. Ljung. Maximum entropy kernels for system identification. *IEEE Transactions on Automatic Control*, 2017.

[19] A. Marconato, M. Schoukens, and J. Schoukens. Filter-based regularisation for impulse response modelling. *IET Control Theory & Applications*, to appear 2016.

[20] K. Zhou, J. C. Doyle, and K. Glover. *Robust and Optimal Control*. Prentice-Hall, Englewood Cliffs, NJ, 1996.

[21] T. Chen and L. Ljung. Implementation of algorithms for tuning parameters in regularized least squares problems in system identification. *Automatica*, 49:2213–2220, 2013.

[22] Nachman Aronszajn. Theory of reproducing kernels. *Transactions of the American mathematical society*, pages 337–404, 1950.

[23] Gianluigi Pillonetto and Alessandro Chiuso. Tuning complexity in regularized kernel-based regression and linear system identification: The robustness of the marginal likelihood estimator. *Automatica*, 58:106–117, 2015.

[24] Claudio Carmeli, Ernesto De Vito, and Alessandro Toigo. Vector valued reproducing kernel hilbert spaces of integrable functions and mercer theorem. *Analysis and Applications*, 4(04):377–408, 2006.

[25] Ingo Steinwart and Andreas Christmann. *Support vector machines*. Springer Science & Business Media, New York, 2008.

[26] R Silverman. Locally stationary random processes. *Information Theory, IRE Transactions on*, 3(3):182–187, 1957.

[27] C. E. Rasmussen and C. K. I. Williams. *Gaussian Processes for Machine Learning*. MIT Press, Cambridge, MA, 2006.

[28] Akira Moiseevich Yaglom. *Correlation theory of stationary and related random functions*. Springer, 1987.

[29] James Mercer. Functions of positive and negative type, and their connection with the theory of integral equations. *Philosophical transactions of the royal society of London. Series A, containing papers of a mathematical or physical character*, pages 415–446, 1909.

[30] Harry Hochstadt. *Integral equations*. John Wiley & Sons, 1973.

[31] Hongwei Sun. Mercer theorem for RKHS on noncompact sets. *Journal of Complexity*, 21(3):337–349, 2005.

[32] Michel Loève. Probability theory, vol. II. *Graduate texts in mathematics*, 1978.

[33] Brett Ninness and Graham C. Goodwin. Estimation of model quality. *Automatica*, 31(12):1771 – 1797, 1995.

[34] L. Ljung. Model validation and model error modeling. In B. Wittenmark and A. Rantzer, editors, *The strm Symposiium on Control*, pages 15 –42, Lund, Sweden, Aug 1999. Studentlitteratur.

[35] G. C. Goodwin, M. Gevers, and B. Ninness. Quantifying the error in estimated transfer functions with application to model order selection. *IEEE Trans. Automatic Control*, 37(7):913–929, 1992.

[36] G. C. Goodwin, J. H. Braslavsky, and M. M. Seron. Non-stationary stochastic embedding for transfer function estimation. *Automatica*, 38:47–62, 2002.

[37] Wolfgang Reinelt, Andrea Garulli, and Lennart Ljung. Comparing different approaches to model error modeling in robust identification. *Automatica*, 38(5):787 – 803, 2002.

[38] *Model Error Modeling and Stochastic Embedding*, volume 48, Beijing, China, 2015.

[39] C. Chen. *Linear system theory and design*. Oxford University Press, New York, 3 edition, 1999.

[40] Torkel Glad and Lennart Ljung. *Control theory: Multivariable and nonlinear methods*. Taylor & Francis, 2000.

[41] T. B. Schön, A. Wills, and B. Ninness. System identification of nonlinear state-space models. *Automatica*, 47(1):39–49, January 2011.

[42] K. J. Åström. *Introduction to stochastic control theory*. Academic Press, New York and London, 1970.

[43] E. T Jaynes. On the rationale of maximum-entropy methods. *Proceedings of the IEEE*, 70(9):939–952, 1982.

[44] F. P. Carli. On the maximum entropy property of the first-order stable spline kernel and its implications. In *IEEE Multi-Conference on Systems and Control*, pages 409–414, Nice/Antibes, France, 2014.

[45] A. Chiuso, T. Chen, L. Ljung, and G. Pillonetto. On the design of multiple kernels for nonparametric linear system identification. In *Proceedings of the IEEE Conference on Decision and Control*, pages 3346–3351, Los Angeles, CA., 2014.

[46] Marc G Genton. Classes of kernels for machine learning: a statistics perspective. *The Journal of Machine Learning Research*, 2:299–312, 2002.

[47] Theodoros Evgeniou, Massimiliano Pontil, and Tomaso Poggio. Regularization networks and support vector machines. *Advances in computational mathematics*, 13(1):1–50, 2000.

[48] F. Cucker and S. Smale. On the mathematical foundations of learning. *American Mathematical Society*, 39(1):1–49, 2002.

[49] T. M. Cover and J. A. Thomas. *Elements of information theory*. John Wiley & Sons, 2012.

[50] A. P. Dempster. Covariance selection. *Biometrics*, 28(1):157–175, 1972.

## Proof of Proposition 3.1

*Part (a)*. For any $n \in \mathbb{N}$ and for any $t_i \in X$, $i = 1, \ldots, n$, let $\bar{b}$ be the column vector containing $b(t_1), \ldots, b(t_n)$ in order. Then the kernel matrix $K^d$ defined by $K_{i,j}^d = k^d(t_i, t_j)$ is rewritten as $K^d = \bar{b}\bar{b}^T$. Clearly, $K^d$ is positive semidefinite, and moreover rank$(K^d) = 1$. So $k^d(t, s)$ is a rank-1 kernel by definition. The identity (14), *Part (b)* and *Part (c)* can be verified in a straightforward way.

## Proof of Proposition 3.2

We give the proof for the CT case with $X = \{t | t \geq 0\}$. The results for the DT case with $X = \mathbb{N}$ hold by noting that the DT kernel is the CT kernel restricted to $\mathbb{N}$.

From (16) and (17), we see the proof will be done if it can be shown that the kernels $k^c(t - s)$ in (16) and (17) are stationary kernels. To show this, note from e.g. [27, page 85] that $e^{-\beta|t-s|}$ with $\beta > 0$ is called the exponential kernel and belongs to the class of Matérn covariance functions with order 1/2. So it remains to show $k^c(t - s)$ in (16) is also a kernel. Note that the exponential kernel is an isotropic stationary (IS) kernel (see Section 3.2.1 for its definition). Below we show that $k^c(t - s)$ in (16) is actually also an IS kernel. As shown in [46], the spectral

representation of an IS kernel $k^c(t,s)$ takes the form of

$$k^c(t,s) = \int_0^\infty \cos(\omega|t-s|)F(d\omega) \qquad \text{(A.1)}$$

where $F$ is any nondecreasing bounded function. For the exponential kernel $e^{-\beta|t-s|}$ with $\beta > 0$, (A.1) is satisfied with the spectral density

$$f(\omega) = \frac{\beta}{\pi(\beta^2 + \omega^2)}.$$

So $k^c(t-s)$ in (16) with $\lambda = \exp(-\beta/2)$ satisfies (A.1) with a well defined spectral density described by

$$\frac{3}{2}\frac{\frac{1}{2}\beta}{\pi((\frac{\beta}{2})^2 + \omega^2)} - \frac{1}{2}\frac{\frac{3}{2}\beta}{\pi((\frac{3\beta}{2})^2 + \omega^2)}$$
$$= \frac{3\beta}{4\pi}\left(\frac{1}{(\frac{\beta}{2})^2 + \omega^2} - \frac{1}{(\frac{3\beta}{2})^2 + \omega^2}\right) \geq 0, \quad \forall \omega$$

Therefore, $k^c(t-s)$ in (16) is an IS kernel. This completes the proof.

**Proof of Proposition 3.3**

We give the proof for the CT case with $X = \{t | t \geq 0\}$. The proof for the DT case with $X = \mathbb{N}$ can be derived in a similar way.

*Part (a)*. Assume that $b(t) \in \mathcal{L}_1$, i.e., $\int_0^\infty b(t)dt < \infty$. Since $|k^c(t-s)| \leq 1$ for any $t, s \geq 0$. Then from

$$\int_0^\infty \left|\int_0^\infty k(t,s)dt\right| ds \leq \int_0^\infty b(s)ds \int_0^\infty b(t)dt < \infty,$$

and by Corollary 2.1, the AMLS kernel (13) is stable.

*Part (b)*. Assume that the AMLS kernel (13) is stable, i.e., $\mathcal{H}_k \subset \mathcal{L}_1$. We need a lemma to prove the result.

**Lemma A.1** [47, p. 16], [48, p. 37] *Assume that $\bar{k}(t,s)$ with $t, s \in X$ is a positive semidefinite kernel and moreover, there exists a sequence of positive numbers $\nu_i$ and linearly independent functions $\psi_i(t)$ such that $\bar{k}(t,s) = \sum_{i=1}^\infty \nu_i \psi_i(t)\psi_i(s)$, where the convergence is absolute and uniform on $Y_1 \times Y_2$ with $Y_1, Y_2$ being any compact subsets of $X$. Then*

$$\mathcal{H}_{\bar{k}} = \left\{g | g = \sum_{i=1}^\infty \mu_i \psi_i, \sum_{i=1}^\infty \frac{\mu_i^2}{\nu_i} < \infty\right\}$$

*and for any $f_1, f_2 \in \mathcal{H}_{\bar{k}}$ with $f_1 = \sum_{i=1}^\infty c_i \psi_i$ and $f_2 = \sum_{i=1}^\infty d_i \psi_i$, the inner product over $\mathcal{H}_{\bar{k}}$ is defined as*

$$\langle f, g\rangle_{\mathcal{H}_{\bar{k}}} = \sum_{i=1}^\infty \frac{c_i d_i}{\nu_i}.$$

Noting (19) and by Lemma A.1, we have

$$\mathcal{H}_{k^c} = \left\{h | h = \sum_{i=1}^\infty \mu_i \phi_i, \sum_{i=1}^\infty \frac{\mu_i^2}{\lambda_i} < \infty\right\}. \qquad \text{(A.2)}$$

Then from (19), the AMLS kernel (13) is rewritten as

$$k(t,s) = \sum_{i=1}^\infty \lambda_i \rho_i(t)\rho_i(s), \rho_i(t) = b(t)\phi_i(t), \qquad \text{(A.3)}$$

and by Lemma A.1, we have

$$\mathcal{H}_k = \left\{g | g = \sum_{i=1}^\infty \mu_i \rho_i, \sum_{i=1}^\infty \frac{\mu_i^2}{\lambda_i} < \infty\right\}. \qquad \text{(A.4)}$$

By (A.4) and (A.2), we have

$$\mathcal{H}_k = \{g | g(t) = b(t)h(t), t \geq 0, h \in \mathcal{H}_{k^c}\}. \qquad \text{(A.5)}$$

Since $\mathcal{H}_k \subset \mathcal{L}_1$, for any $g \in \mathcal{H}_k$, $g \in \mathcal{L}_1$, i.e.,

$$\int_0^\infty |g(t)|dt = \int_0^\infty b(t)|h(t)|dt < \infty, \qquad \text{(A.6)}$$

where the equality follows from (A.5) and $h \in \mathcal{H}_{k^c}$ is the function associated with $g \in \mathcal{H}_k$.

Since $k^c$ is not stable, $\mathcal{H}_{k^c} \not\subset \mathcal{L}^1$ and there exists $h \in \mathcal{H}_{k^c}$ but $h \notin \mathcal{L}^1$. For such $h$, if there exits an $\epsilon > 0$ such that $b(t) \geq \epsilon$ for all $t \geq 0$, then

$$\int_0^\infty |g(t)|dt = \int_0^\infty b(t)|h(t)|dt \geq \epsilon \int_0^\infty |h(t)|dt = \infty,$$

which contradicts with (A.6). This completes the proof.

**Proof of Proposition 4.1**

To prove the result, we need a lemma.

**Lemma A.2** [42,40] *Consider a zero mean stationary Gaussian process $h(t)$ with covariance function $k^c(t-s)$. Assume that $k^c(t-s)$ has rational power spectral density $\Psi(\omega)$[8]. Then the following results hold:*

- *For the DT case, there exists a rational function $G$ which has all poles inside the unit circle and all zeros inside or on the unit circle such that*

$$\Psi(\omega) = G(e^{i\omega})G(e^{-i\omega}) \qquad \text{(A.7)}$$

---

[8] Check the statement of Theorem 4.1 for the definition of rational power spectral density.

where $|\cdot|$ denotes the modulus of a complex number. Moreover, let the quintuple $(\bar{A}, \bar{B}, \bar{C}, \bar{D}, \bar{Q})$ be a state space realization of $G(e^{i\omega})$, i.e., $G(e^{i\omega}) = \bar{C}(e^{i\omega}I_{\dim(\bar{A})} - \bar{A})^{-1}\bar{B} + \bar{D}$, where $\dim(\bar{A})$ is the dimension of $\bar{A}$ and $I_{\dim(\bar{A})}$ is the identity matrix with dimension $\dim(\bar{A})$. Then the stationary Gaussian process $h(t)$ has the following state space representation

$$x(t+1) = \bar{A}x(t) + \bar{B}w(t), x(0) \sim \mathcal{N}(0, \bar{Q}) \quad \text{(A.8a)}$$
$$h(t) = \bar{C}x(t) + \bar{D}w(t) \quad \text{(A.8b)}$$

where $w(t)$ is the zero mean white Gaussian noise with unit variance and $\bar{Q}$ is the solution of the Lyapunov equation

$$\bar{Q} = \bar{A}\bar{Q}\bar{A}^T + \bar{B}\bar{B}^T \quad \text{(A.8c)}$$

- For the CT case, there exists a rational function $G$ which has all poles in the left half plane and all zeros in the left half plane or on the imaginary axis such that

$$\Psi(\omega) = G(i\omega)G(-i\omega). \quad \text{(A.9)}$$

Moreover, let the quintuple $(\bar{A}, \bar{B}, \bar{C}, \bar{D}, \bar{Q})$ be a state space realization of $G(i\omega)$, i.e., $G(i\omega) = \bar{C}(i\omega I_{\dim \bar{A}} - \bar{A})^{-1}\bar{B} + \bar{D}$. Then the stationary Gaussian process $h(t)$ has the following state space representation

$$\dot{x}(t) = \bar{A}x(t) + \bar{B}w(t), x(0) \sim \mathcal{N}(0, \bar{Q}) \quad \text{(A.10a)}$$
$$h(t) = \bar{C}x(t) + \bar{D}w(t) \quad \text{(A.10b)}$$

where $w(t)$ is the zero mean white Gaussian noise with unit power spectral density and $\bar{Q}$ is the solution of the Lyapunov equation

$$\bar{A}\bar{Q} + \bar{Q}\bar{A}^T + \bar{B}\bar{B}^T = 0 \quad \text{(A.10c)}$$

**Proof:** For the DT case, the first part is a result of the Spectral Factorization Theorem [42, Theorem 3.1 in Chapter 4] and the second part is a result of the Representation Theorem [42, Theorem 3.2 in Chapter 4] and [40, Theorem 5.3 and Equation (5.92)]. For the CT case, the first part is a result of the Spectral Factorization Theorem [42, Theorem 5.1 in Chapter 4] and the second part is a result of the Representation Theorem [42, Theorem 5.2 in Chapter 4] and [40, Theorem 5.3]. $\diamond$.

For the DT case, let $z(t) = c^{\frac{1}{2}}\lambda^t x(t)$. Then simple calculation shows that the AMLS kernel (13) with $k^d(t, s) = c\lambda^{t+s}$ is in the form of (24) with $A = \lambda\bar{A}$, $B = \lambda\bar{B}$, $C = \bar{C}$, $D = \bar{D}$, $Q = c\bar{Q}$ and $b(t) = c^{\frac{1}{2}}\lambda^t$. For the CT case, let $z(t) = c^{\frac{1}{2}}e^{-\beta t}x(t)$ with $\lambda = e^{-\beta}$. Then simple calculation shows that the AMLS kernel (13) with $k^d(t, s) = c\lambda^{t+s}$ is in the form of (24) with $A = \bar{A} - \beta I_{\dim \bar{A}}$, $B = \bar{B}$, $C = \bar{C}$, $D = \bar{D}$, $Q = c\bar{Q}$ and $b(t) = c^{\frac{1}{2}}e^{-\beta t}$.

**Proof of Proposition 4.2**

We only sketch the proof for the SS kernel and the proof for the DC kernel can be derived in a similar way.

For the DT case, the IS kernel $k^c(t - s)$ in (16) has the power spectral density

$$\begin{aligned}
\Psi(\omega) &= \sum_{\tau=-\infty}^{+\infty} k^c(\tau)e^{-i\omega\tau} \\
&= \frac{3}{2}\frac{1-\lambda^2}{(1-\lambda e^{-i\omega})(1-\lambda e^{i\omega})} - \frac{1}{2}\frac{1-\lambda^6}{(1-\lambda^3 e^{-i\omega})(1-\lambda^3 e^{i\omega})} \\
&= \frac{(1-\lambda^2)^3}{2}\frac{2+2\lambda^2+\lambda(e^{-i\omega}+e^{-i\omega})}{(1-\lambda e^{-i\omega})(1-\lambda e^{i\omega})(1-\lambda^3 e^{-i\omega})(1-\lambda^3 e^{i\omega})} \\
&= \frac{(1-\lambda^2)^3}{2}\frac{(\bar{a}e^{-i\omega}+\bar{b})(\bar{a}e^{i\omega}+\bar{b})}{(1-\lambda e^{-i\omega})(1-\lambda e^{i\omega})(1-\lambda^3 e^{-i\omega})(1-\lambda^3 e^{i\omega})}
\end{aligned}$$

By spectral factorization technique, we have $\Psi(\omega) = G(e^{i\omega})\overline{G(e^{i\omega})}$ with

$$G(e^{i\omega}) = \sqrt{\frac{(1-\lambda^2)^3}{2}}\frac{(\bar{a}e^{-i\omega}+\bar{b})}{(1-\lambda e^{-i\omega})(1-\lambda^3 e^{-i\omega})}.$$

For the CT case, the IS kernel $k^c(t - s)$ has the power spectral density

$$\begin{aligned}
\Psi(\omega) &= \int_{-\infty}^{+\infty}\left(\frac{3}{2}e^{-\frac{1}{2}\beta|\tau|} - \frac{1}{2}e^{-\frac{3}{2}\beta|\tau|}\right)e^{-iw\tau}d\tau \\
&= \left|\frac{3^{\frac{1}{2}}\beta^{\frac{3}{2}}}{(i\omega+\frac{1}{2}\beta)(i\omega+\frac{3}{2}\beta)}\right|^2 \triangleq |G(i\omega)|^2.
\end{aligned}$$

From the realization theory of linear systems, see e.g., [39], we can derive the corresponding state-space model representation of the IS kernel $k^c(t - s)$, based on which and by using the argument in the end of the proof of Proposition 4.1, we derive (27) and (29), respectively.

**Proof of Proposition 4.3**

For the DT case, note that there exists $l > 0$ such that

$$\begin{aligned}
&|CA^tQ(A^s)^TC^T| \le l|\bar{\lambda}|^{t+s}, \\
&|D^2b(t)b(s)\delta(t-s)| \le l|\bar{\lambda}|^{\frac{t+s}{2}}, \\
&\left|Db(t)\sum_{k=0}^{s-1}\delta(t-k)b(k)B^T(A^{s-1-k})^TC^T\right| \le l|\bar{\lambda}|^{\frac{t+s}{2}}, \\
&\left|\sum_{k=0}^{\min\{t,s\}-1}b(k)^2CA^{t-1-k}BB^T(A^{s-1-k})^TC^T\right| \\
&\quad\le l(|\bar{\lambda}|^{t+s} + |\bar{\lambda}|^{t+s-\min\{t,s\}}).
\end{aligned}$$

17

Then it is easy to show $\sum_{t=1}^{\infty}\sum_{s=1}^{\infty}\left|k^{\mathrm{SI}}(t,s)\right| < \infty$, and thus the SI kernel (25a) is stable by Corollary 2.1.

For the CT case, note that there exists $l > 0$ such that

$$|Ce^{At}Q(e^{As})^T C^T| \le le^{-|\mathrm{Re}(\bar{\lambda})|(t+s)},$$

$$\left| C\int_0^{\min\{t,s\}} b^2(\tau)e^{A(t-\tau)}BB^T(e^{A(s-\tau)})^T d\tau C^T \right|$$
$$\le l(e^{-|\mathrm{Re}(\bar{\lambda})|(t+s)} + e^{-|\mathrm{Re}(\bar{\lambda})|(t+s-\min\{t,s\})}).$$

Then it is easy to show $\int_0^{\infty}\int_0^{\infty}\left|k^{\mathrm{SI}}(t,s)\right| dtds < \infty$, and thus the SI kernel (25b) is stable by Corollary 2.1.

**Proof of Proposition 4.4**

We only give the proof for the case $D \ne 0$ and the proof for the case (b) is similar and thus omitted.

By chain rule, the differential entropy in (33) becomes

$$H(\bar{z}(0),\bar{g}(0),\bar{g}(1),\cdots,\bar{g}(s)) = H(\bar{z}(0)) + H(\bar{g}(0)|\bar{z}(0))$$
$$+\sum_{t=1}^{s} H(\bar{g}(t)|\bar{z}(0),\bar{g}(0),\cdots,\bar{g}(t-1)) \qquad (A.11)$$

Note from (32) that

$$H(\bar{g}(0)|\bar{z}(0)) = H(f(0)Db(0)+C\bar{z}(0)|\bar{z}(0))$$
$$= H(f(0)Db(0)|\bar{z}(0))$$
$$= H(f(0)|\bar{z}(0)) + \log|Db(0)|$$

where the first equation follows because translation does not change the differential entropy and the second equation follows because of the scaling property of differential entropy. Analogously, we have

$$H(\bar{g}(t)|\bar{z}(0),\bar{g}(0),\cdots,\bar{g}(t-1))$$
$$= H(\bar{g}(t)|\bar{z}(0),f(0),\cdots,f(t-1))$$
$$= H(f(t)Db(t)|\bar{z}(0),f(0),\cdots,f(t-1))$$
$$= H(f(t)|\bar{z}(0),f(0),\cdots,f(t-1)) + \log|Db(t)|$$

Therefore, (A.11) is rewritten as

$$H(\bar{z}(0),\bar{g}(0),\bar{g}(1),\cdots,\bar{g}(s)) = H(\bar{z}(0)) + H(f(0)|\bar{z}(0))$$
$$+\sum_{t=1}^{s} H(f(t)|\bar{z}(0),f(0),\cdots,f(t-1)) + \sum_{t=0}^{s} \log|Db(t)|$$
$$= H(\bar{z}(0),f(0),f(1),\cdots,f(s)) + \sum_{t=0}^{s} \log|Db(t)|$$

Since the second term in the above equation is independent of $\bar{z}(0)$ and $f(t)$, $t = 0,\cdots,s$, the MaxEnt problem

(33) is equivalently rewritten as

$$\underset{\bar{z}(0),f(t)}{\mathrm{maximize}} \quad H(\bar{z}(0),f(0),f(1),\cdots,f(s)) \qquad (A.12a)$$
subject to
$$\mathbb{E}(\bar{z}(0)) = 0, \mathbb{E}(f(t)) = 0,$$
$$\mathbb{C}ov(\bar{z}(0)) = Q, \mathbb{V}(f(t)) = 1, t = 0,\cdots,s. \quad (A.12b)$$

Note that the constraints in (A.12) are separable with respect to $\bar{z}(0), f(t)$, $t = 0,\cdots,s$, and that $H(\bar{z}(0),f(0),f(1),\cdots,f(s))$ is maximized if $\bar{z}(0), f(0), f(1),\cdots,f(s)$ are independent with each other. Therefore, (A.12) is equivalent to

$$\underset{\bar{z}(0)}{\mathrm{maximize}} \quad H(\bar{z}(0))$$
$$\text{subject to } \mathbb{E}(\bar{z}(0)) = 0, \mathbb{C}ov(\bar{z}(0)) = Q$$

and for $t = 0,\cdots,s$,

$$\underset{f(t)}{\mathrm{maximize}}\, H(f(t))$$
$$\text{subject to } \mathbb{E}(f(t)) = 0, \mathbb{V}(f(t)) = 1$$

It is well known that multivariate normal distribution maximizes the differential entropy over all distributions with the same covariance, see e.g., [49, Theorem 8.6.5]. Then we have the optimal solution to (A.12) is that $\bar{z}(0) \sim \mathcal{N}(0,Q)$ and $f(t) \sim \mathcal{N}(0,1)$, $t = 0,\cdots,s$, and moreover, $\bar{z}(0), f(0), f(1),\cdots,f(s)$ are independent with each other. Finally, (32) implies $\bar{g}(t) = CA^t\bar{z}(0)+\sum_{k=0}^{t-1} CA^{t-1-k}Bb(k)f(k)+Db(t)f(t)$, same as $g(t)$ in (25a). This completes the proof.

**Proof of Corollary 4.1**

From Proposition 4.2, the DC kernel is a SI kernel with $A,B,C,D,Q,b(t)$ given in (28), which implies that Theorem 4.4 holds for the DC kernel. Comparing (32) with (36) shows that the proof is completed if we can show

$$\mathbb{C}ov(\bar{z}(0)) = Q \quad \Longrightarrow \quad \mathbb{V}(\bar{g}(0)) = c,$$
$$\mathbb{V}(f(t)) = 1 \Longrightarrow \mathbb{V}(\bar{g}(t) - \lambda^{1/2}\rho\bar{g}(t-1))$$
$$= c(1 - \rho^2)\lambda^t, t = 1,\cdots,s \quad (A.13)$$

The first line of (A.13) is straightforward. In what follows, we consider the second line of (A.13). Clearly, the results holds for the case either $\lambda = 0$ or $\rho = 0$. So we further consider below the case where $\lambda \ne 0$ and $\rho \ne 0$.

It is easy to see that for $t = 0, \ldots, s-1$,

$$\bar{g}(t+1) = ACA^t\bar{z}(0) + A\sum_{k=0}^{t-1}CA^{t-1-k}Bb(k)f(k)$$

$$+ ACA^{-1}Bb(t)f(t) + Db(t+1)f(t+1)$$

$$= ACA^t\bar{z}(0) + A\sum_{k=0}^{t-1}CA^{t-1-k}Bb(k)f(k)$$

$$+ ADb(t)f(t) + Db(t+1)f(t+1)$$

$$= A\bar{g}(t) + Db(t+1)f(t+1)$$

where the second equality holds because $CA^{-1}B = D$. Therefore (A.13) holds, which completes the proof.

**Proof of Proposition 4.5**

We need a lemma to prove this result.

**Lemma A.3** *Consider the pth order Gaussian Markov process (38). Then the following results hold:*

*(a) For any $t \in \mathbb{N}$ and $k > p$, $x(t)$ and $x(t+k)$ are conditionally independent given $x(s)$ with $s \neq t$, $s \neq t+k$ and $s \in \mathbb{N}$.*

*(b) Let $i_1 < i_2 < \cdots < i_n$ with $n > p+1$ be any strictly increasing subsequence taken from $\mathbb{N}$. Then the covariance matrix of $[x(i_1) \ \cdots \ x(i_n)]^T$ is a p-band matrix.*

**Proof:** *Part (a) follows from the pth order Markov property of $x(t)$ and the fact that $w(t)$ is white Gaussian noise. To prove Part (b), recall from e.g., [50] that, for Gaussian random variables, the zeros in the inverse of the covariance matrix indicate conditional independence of the two corresponding elements conditioned on the remaining ones. To be specific, consider a Gaussian random variable with mean $m$ and covariance matrix $K$:*

$$\begin{bmatrix} x \ y \ z \end{bmatrix}^T \sim \mathcal{N}(m, K), \quad \text{(A.14)}$$

*where $x, y \in \mathbb{R}$, $z \in \mathbb{R}^n$ and $K \in \mathbb{R}^{(n+2)\times(n+2)}$. Then for a given $z$, $x, y$ are conditionally independent if and only if $(K^{-1})_{1,2} = (K^{-1})_{2,1} = 0$, where $(K^{-1})_{i,j}$ denote the $(i,j)$th element of $K^{-1}$. Therefore, Part (b) follows from the above observation and part (a).* ◇

Note from Proposition 4.2 that $g(t)$ can be written in the following form

$$z(t+1) = \lambda^{\frac{1}{2}}\rho z(t) + \lambda^{\frac{1}{2}}c^{\frac{1}{2}}\lambda^{\frac{t}{2}}w(t), z(0) \sim \mathcal{N}(0, \frac{c}{1-\rho^2})$$

$$g(t) = \rho(1-\rho^2)^{\frac{1}{2}}z(t) + (1-\rho^2)^{\frac{1}{2}}c^{\frac{1}{2}}\lambda^{\frac{t}{2}}w(t).$$

which implies that (39) holds. Clearly, (39) is in the form of (38a) with $p = 1$ and thus $g(t)$ is a Markov process

with order 1. Finally, the result that the kernel matrix of the DC kernel has a 1-band matrix inverse follows from Lemma A.3.

**Proof of Proposition 4.6**

By using the realization of $G^0(q)$ in controllable canonical form, the DT SI kernel (24) is written as follows:

$$z(t+1) = \begin{bmatrix} -a_1 \ \cdots \ -a_{n-1} \ -a_n \\ I_n \qquad\qquad \mathbf{0}_{n\times 1} \end{bmatrix} z(t)$$

$$+ \begin{bmatrix} 1 \\ \mathbf{0}_{n\times 1} \end{bmatrix} b(t)w(t)$$

$$g(t) = \begin{bmatrix} \bar{b} \ 0 \ 0 \ \cdots \ 0 \end{bmatrix} z(t), z(0) \sim \mathcal{N}(0, Q).$$

From the above equation, we have

$$g(t+1) = -\bar{b}\sum_{i=1}^{n} a_i z_i(t) + \bar{b}b(t)w(t) \qquad \text{(A.15)}$$

where $z_i(t)$ is the $i$th element of $z(t)$. Note that

$$\bar{b}z_i(t) = g(t-i+1), i = 1, \cdots, n$$

Therefore, (A.15) becomes

$$g(t+1) = -\sum_{i=1}^{n} a_i g(t-i+1) + \bar{b}b(t)w(t)$$

which is in the form of (38b) with order $n$. Thus the Gaussian process $g(t)$ in (24) is a Markov process with order $n$ and the fact that the kernel has an $n$-band matrix inverse follows from Lemma A.3.