

Mathematical foundations of stable RKHSs

Mauro Bisiacco^a

^a*Department of Information Engineering, University of Padova, Padova, Italy (e-mail: bisiacco@dei.unipd.it)*

Gianluigi Pillonetto^b

^b*Department of Information Engineering, University of Padova, Padova, Italy (e-mail: giapi@dei.unipd.it)*

Abstract

Reproducing kernel Hilbert spaces (RKHSs) are key spaces for machine learning that are becoming popular also for linear system identification. In particular, the so-called stable RKHSs can be used to model absolutely summable impulse responses. In combination e.g. with regularized least squares they can then be used to reconstruct dynamic systems from input-output data. In this paper we provide new structural properties of stable RKHSs. The relation between stable kernels and other fundamental classes, like those containing absolutely summable or finite-trace kernels, is elucidated. These insights are then brought into the feature space context. First, it is proved that any stable kernel admits feature maps induced by a basis of orthogonal eigenvectors in ℓ_2 . The exact connection with classical system identification approaches that exploit such kind of functions to model impulse responses is also provided. Then, the necessary and sufficient stability condition for RKHSs designed by formulating kernel eigenvectors and eigenvalues is obtained. Overall, our new results provide novel mathematical foundations of stable RKHSs with impact on stability tests, impulse responses modeling and computational efficiency of regularized schemes for linear system identification.

Key words: linear system identification; BIBO stability; stable reproducing kernel Hilbert spaces; kernel-based regularization; regularized least squares

1 Introduction

Reproducing kernel Hilbert spaces (RKHSs) are particular spaces of functions in one-to-one correspondence with the class of positive semidefinite kernels. While RKHS theory has been mainly developed in the fifties [2,6], such spaces have found first important applications in the eighties in the context of statistics and computer vision [7,54,48]. They were then brought to the attention of machine learning community in [28]. Since then, they have become a fundamental tool for function estimation. Estimators based on RKHSs include smoothing splines [54], regularization networks [48] and support vector machines [25,53]. Combinations with deep networks are also described in [20,5].

The importance of RKHSs for function estimation from sparse and noisy data arises from several facts. First, a RKHS \mathcal{H} inherits its properties from the associated kernel K . This is important for modeling purposes since all the expected function properties can be encoded in the kernel design. For instance, a regular kernel induces an RKHS of continuous functions whose norm can be used as regularizer to penalize solutions with unphysical oscillations. Indeed, the most important kernel-based estimators optimize objectives containing a loss that accounts for adherence to experimental data and the RKHS norm that restores well-posedness. Another fundamental aspect is that the kernel can include in an implicit way a very large (possibly infinite) number of basis functions, leading to very flexible and computable models. This result has also connection with the following important mathematical fact. Let \mathcal{X} be the regressor space, i.e. the domain of the functions $f : \mathcal{X} \rightarrow \mathbb{R}$ contained in \mathcal{H} . Then, given any positive semidefinite kernel K , there always exists at least one inner-product space \mathcal{F} and one feature

map $\phi : \mathcal{X} \rightarrow \mathcal{F}$ such that¹

$$K(x, y) = \langle \phi(x), \phi(y) \rangle_{\mathcal{F}}, \quad \phi : \mathcal{X} \rightarrow \mathcal{F}. \quad (1)$$

Above, the components of $\phi(x)$ are the basis functions induced by the kernel, e.g. $\phi_1(x) = 1, \phi_2(x) = x, \phi_3(x) = x^2, \dots$ describe polynomial models. Kernels can thus define expressive spaces by (implicitly) mapping the space of the regressors \mathcal{X} into high-dimensional feature spaces where linear machines can be employed. Nonlinear algorithms can be reduced to linear ones without even knowing explicitly the feature map. In fact, under mild assumptions, the kernel-based estimate is given by the sum of a finite number of kernel sections $K(x, \cdot)$ centred on the observed regressors [54,51].

Control community's interest has been recently addressed to RKHSs tailored for linear system identification. While function models adopted in machine learning typically embed information e.g. on regularity, periodicity or sparsity, the new spaces hinge on kernels accounting for dynamic systems features. Examples are the so-called stable spline, TC and DC kernels that incorporate exponential stability [44,45,17,19], see also [12,13,38,39] for even more sophisticated models. They define regularized least squares schemes that can outperform conventional parametric identification [46,35,9,16,42]. All of these kernels belong to the more general class of (BIBO) stable kernels that induce RKHSs of absolutely summable impulse responses. A fundamental characterization of these kernels has been known in the literature at least since 2006 [11]. It says that a kernel K is stable if and only if it induces a bounded integral operator mapping the space ℓ_∞ of essentially bounded functions into the space ℓ_1 of absolutely summable functions, see also [24,18]. This result is the starting point of this paper. Building upon it, new structural properties of stable RKHSs will be obtained working in discrete-time (\mathcal{X} becomes the set of natural numbers). In particular, the main contributions of the paper are the following ones.

First, we will obtain fundamental RKHSs inclusion properties that shed new light on the relationships between stable kernels and e.g. absolutely summable and finite-trace kernels. This result defines in a natural and simple way new stability tests on several classes of kernels. It also contains, as immediate corollaries, some instability results obtained in the literature through ad-hoc theorems, e.g. regarding the translation invariant class that contains the popular Gaussian kernel [36,24,40].

As for the second contribution, let ℓ_2 be the space of squared summable sequences with the usual inner product. Then, we show that any stable kernel admits a *spectral (Mercer) feature map* $\phi : \mathcal{X} \rightarrow \ell_2$ with

$$\phi(x) = \{ \sqrt{\lambda_i} \rho_i(x) \}_{i=1}^{\infty}$$

where λ_i and ρ_i are, respectively, the kernel eigenvalues and eigenfunctions forming an orthonormal basis in ℓ_2 . The fact that any stable RKHS is generated by such basis provides the fundamental link with the important literature on impulse response estimation via orthonormal functions [55,56,37,30]. Furthermore, under an algorithmic viewpoint, many efficient machine learning procedures exploit truncated Mercer expansions for approximating the kernel, e.g. see [57,32,29] and also [59,47] for discussions on their optimality in a stochastic framework. These works trace back to the so-called Nyström method where an integral equation is replaced by finite-dimensional approximations [3,4]. For system identification, the works [41,10] have shown that a relatively small number of eigenfunctions (w.r.t. the data set size) can capture impulse responses regularized estimates. Our result shows that any stable kernel is amenable to these fast computational schemes. To exploit them, closed-form expressions of the λ_i and ρ_i are desirable. Determining the spectrum of K is however far from trivial in general but numerical approximations can be adopted. In this regard, we show that singular value decompositions applied to truncated stable kernels generate a sequence of spectra convergent in ℓ_2 to the correct one. This can be seen as a novel convergence result for a Nyström-type method on unbounded kernel domains.

Third, having established that any stable RKHS is generated by a basis of ℓ_2 , the question is however which kind of orthonormal functions and of their combinations lead to stable kernels. This motivates the study of stability conditions for models built through feature maps and (1). This route loses the advantage of implicit encoding since it rarely leads to closed-form kernel expressions (this is the dual problem of the Mercer expansion). However, such issue is very relevant. In fact, recent literature has shown that important dynamic system features, like the presence of resonances, can be conveniently described using feature maps e.g. induced by Kautz models [15,23]. Then, we will provide the necessary and sufficient stability condition for kernels defined by Mercer expansions. This new outcome should be taken into account when formulating any linear system model whose aim is to combine orthogonal basis functions in ℓ_2 with stability information via kernel-based regularization.

¹ One explicit example is the RKHS map $\phi_{\mathcal{H}} : \mathcal{X} \rightarrow \mathcal{H}$ such that $\phi_{\mathcal{H}}(x) = K(x, \cdot)$. It always satisfies (1) in view of the so-called reproducing property [2].

So, overall, our results have impact on stability tests, impulse responses modeling and computational efficiency issues. To illustrate them, the paper is organized as follows. Section 2 reports a brief overview on (stable) RKHSs setting up also some notation. In Section 3 we obtain new inclusion properties of some notable kernels classes that provide fundamental insights on the structure of stable kernels. Section 4 shows that any stable kernel admits a Mercer expansion in ℓ_2 and discusses the link with impulse response estimation via orthonormal bases. It also shows how to numerically recover kernel eigenfunctions and eigenvalues. In Section 5 the necessary and sufficient condition for RKHS stability in the Mercer feature space is worked out. Conclusions then end the paper while the proof of all the new theorems are gathered in Appendix.

2 Overview on RKHSs and stability condition

We are interested in spaces of functions containing discrete-time impulse responses of causal systems. Hence, the function domain is the set of natural numbers \mathbb{N} . In addition, the elements of the space can be also seen as sequences containing impulse response coefficients.

We will consider in particular the so-called Reproducing Kernel Hilbert Spaces (RKHSs). They are in one-to-one correspondence with positive semidefinite kernels that, in our setting, map $\mathbb{N} \times \mathbb{N}$ into the real line. However, in view of the nature of the domain, in what follows it is more convenient to see the kernel as an infinite-dimensional matrix with the (i, j) -entries denoted by K_{ij} . The positive semidefinite constraints then imply that, for any choice of integers $\{p_1, \dots, p_m\}$, the $m \times m$ matrix A , with $A_{ij} = K_{p_i p_j}$, is symmetric and positive semidefinite.

As already recalled, the RKHS inherits the properties of a kernel. Indeed, the values K_{ij} can be interpreted as a similarity measure between the i -th and the j -th element of the sequence. In linear system identification, the interest is in particular addressed to (BIBO) stable kernels. They induce RKHSs containing only absolutely summable vectors. To introduce them, let ℓ_∞ and ℓ_1 be the spaces of bounded and absolutely summable sequences of real numbers, respectively, i.e.

$$\ell_\infty = \left\{ \{u_i\}_{i \in \mathbb{N}} \text{ s.t. } \|u\|_\infty < \infty \right\},$$

and

$$\ell_1 = \left\{ \{u_i\}_{i \in \mathbb{N}} \text{ s.t. } \|u\|_1 < \infty \right\},$$

with

$$\|u\|_\infty = \sup_{i \in \mathbb{N}} |u_i| \quad \text{and} \quad \|u\|_1 = \sum_{i \in \mathbb{N}} |u_i|.$$

Now, note also that the kernel K defines an acausal linear time-varying system, often called kernel operator in the literature: given an input (sequence) u , the output at instant i is $\sum_{j=1}^{\infty} K_{ij} u_j$. Then, using notation of ordinary algebra to handle infinite-dimensional objects, the output can be indicated by Ku with u an infinite-dimensional (column) vector. With this in mind, the following fundamental theorem reports the necessary and sufficient condition for RKHS stability.

Theorem 1 (RKHS stability [11]) *Let \mathcal{H} be the RKHS induced by K . Then, it holds that*

$$\mathcal{H} \subset \ell_1 \iff Ku \in \ell_1 \quad \forall u \in \ell_\infty. \tag{2}$$

■

The stability condition is equivalent to requiring that the kernel operator is a bounded (continuous) map between ℓ_∞ and ℓ_1 , see [8].

Remark 2 *The third fundamental space used in this paper, already mentioned in the introduction, is that containing squared summable sequences, i.e.*

$$\ell_2 = \left\{ \{u_i\}_{i \in \mathbb{N}} \text{ s.t. } \|u\|_2 < \infty \right\},$$

with

$$\|u\|_2^2 = \sum_{i \in \mathbb{N}} u_i^2.$$

In particular, two types of kernel operators induced by K will be encountered during our analysis. The first one is that described above mapping ℓ_∞ into ℓ_1 while the second one is that mapping ℓ_2 into ℓ_2 itself.

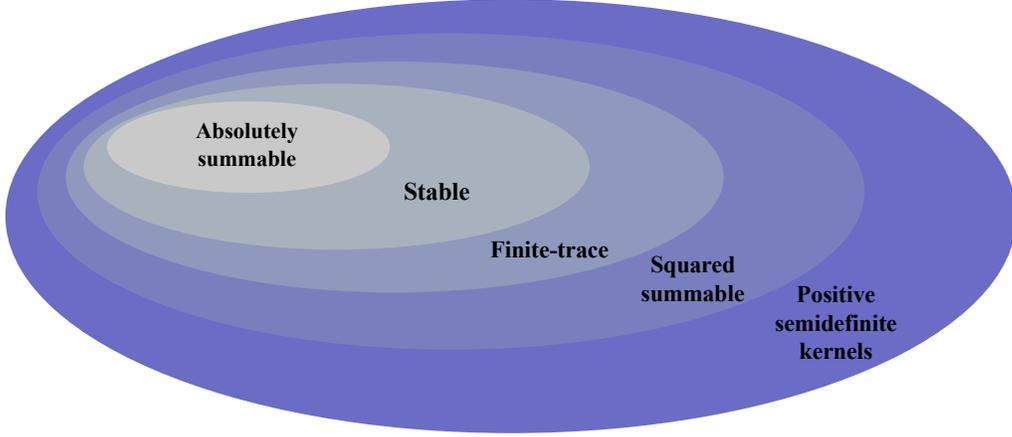


Fig. 1. Inclusion properties of different kernel classes

3 Inclusion properties of some notable kernels classes

In this section we derive new relationships between stable kernels and other fundamental classes. Let \mathcal{S}_s be the set containing all the stable RKHSs. Then, we also consider

- the set \mathcal{S}_1 containing all the RKHSs induced by absolutely summable kernels, i.e. satisfying the constraint

$$\sum_{ij} |K_{ij}| < +\infty;$$

- the set \mathcal{S}_{ft} of RKHSs associated to finite-trace kernels that are characterized by

$$\sum_i K_{ii} < +\infty;$$

- the set \mathcal{S}_2 induced by squared summable kernels, i.e. satisfying

$$\sum_{ij} K_{ij}^2 < +\infty.$$

The following result then holds.

Theorem 3 *One has*

$$\mathcal{S}_1 \subset \mathcal{S}_s \subset \mathcal{S}_{ft} \subset \mathcal{S}_2 \tag{3}$$

■

Fig. 1 provides a graphical description of Theorem 3 in terms of inclusions of kernels classes. Some comments about its meaning, under the perspective of stability tests, are now in order.

Regarding $\mathcal{S}_1 \subset \mathcal{S}_s$, while it is trivial to show that kernel absolute summability implies stability, the notable fact is that such inclusion is strict. In this regard, recall that in [11,24], immediately after reporting Theorem 1, the authors mentioned kernel absolute summability as a sufficient condition, with the desire to formulate a (in some sense) simpler stability test. Its necessity was however left as an open problem. Many papers have then cited and exploited kernel summability as a stability check without answering such question, e.g. see [22,14,27,12]. Theorem 3 points out that the equivalence does not hold. So, one cannot conclude that a kernel is unstable from the sole failure of absolute summability.

The relation $\mathcal{S}_s \subset \mathcal{S}_{ft}$ means that the set of finite-trace kernels contains the stable class. Also this inclusion is strict, hence the analysis of the trace is useful only to prove that a given RKHS is not contained in ℓ_1 . This has however interesting consequences. For instance, in [36] the instability of the Gaussian kernel

$$K_{ij} = e^{-(i-j)^2}$$

was proved exploiting a complex RKHS representation through a generalization of the Weyl inner product for the homogeneous polynomial space. In [24][Appendix] this result was greatly extended by proving that all the RKHSs induced by translation invariant kernels, i.e. of the form

$$K_{ij} = h(i-j)$$

(with h satisfying the positive semidefinite constraints) are not contained in ℓ_1 . The Schoenberg representation theorem was used, see p. 3309 of [24]. All of these ad-hoc theorems are now trivial corollaries of Theorem 3 since the trace of a translation invariant kernel is $\sum_i K_{ii} = \sum_i h(0)$ and it always diverges unless h is the null function. Even more importantly, many other instability results become immediately available. One can e.g. claim that all the kernels whose diagonal elements satisfy $K_{ii} \propto i^{-\delta}$ are unstable if $\delta \leq 1$.

Finally, the strict inclusion $\mathcal{S}_{ft} \subset \mathcal{S}_2$ shows that a stability check relying on kernel squared summability does not make much sense. In fact, the finite-trace test is in any case both more powerful and simpler to perform.

4 Mercer expansions of stable kernels

4.1 Mercer feature maps for stable kernels

Recalling also Remark 2, we are now interested in the operator induced by a stable kernel as a map from ℓ_2 into itself. A kernel operator is compact if it maps any bounded sequence $\{v_i\}$ into a sequence $\{Kv_i\}$ from which a convergent subsequence can be extracted [50,58]. Theorem 3 ensures that any stable kernel K is finite-trace. Combining this fact with Lemma 16 present in Appendix, one obtains the following result.

Theorem 4 *Any operator induced by a stable kernel is self-adjoint, positive semidefinite and compact as a map from ℓ_2 into ℓ_2 itself.* ■

The above theorem is important because it allows us to associate to any stable kernel a Mercer feature map built through an orthonormal basis of ℓ_2 . In particular, the following result is a direct consequence of the spectral theorem [26] (applied to kernels defined over $\mathbb{N} \times \mathbb{N}$) that holds indeed by virtue of Theorem 4.

Proposition 5 (Representation of stable kernels) *Let K be stable. Then, there always exists an orthonormal basis of ℓ_2 composed by eigenvectors $\{\rho_i\}$ of K with corresponding eigenvalues $\{\lambda_i\}$, i.e.*

$$K\rho_i = \lambda_i\rho_i, \quad i = 1, 2, \dots$$

In addition, the spectral Mercer feature map $\phi : \mathbb{N} \rightarrow \ell_2$ with

$$\phi(x) = \{\sqrt{\lambda_i}\rho_i(x)\}_{i=1}^{\infty}$$

is always well-defined and each (x, y) entry of K admits the representation

$$K_{xy} = \langle \phi(x), \phi(y) \rangle_2 = \sum_{i=1}^{+\infty} \lambda_i \rho_i(x) \rho_i(y), \quad (4)$$

where $x, y \in \mathbb{N}$. ■

The pointwise convergence of the kernel expansion (4) stated in Proposition 5, combined with the same arguments used in [21][Chapter 3, Theorem 4] or [52][Theorem 1], allows us to obtain the following characterization of any stable RKHS.

Proposition 6 (Representation of stable RKHSs) *Let K be stable and assume that any kernel eigenvalue satisfies $\lambda_i > 0$. Then, the stable RKHS associated to K always admits the representation*

$$\mathcal{H} = \left\{ f = \sum_{i=1}^{\infty} a_i \rho_i \text{ s.t. } \sum_{i=1}^{\infty} \frac{a_i^2}{\lambda_i} < +\infty \right\}, \quad (5)$$

where the ρ_i are the eigenvectors of K forming an orthonormal basis of ℓ_2 . ■

Remark 7 *In Proposition 6 we have assumed that all the eigenvalues of the stable kernel K are strictly positive so that \mathcal{H} is infinite-dimensional. If some eigenvalue is null, \mathcal{H} is spanned only by the eigenvectors associated to non-null λ_i . If only a finite number of λ_i is different from zero, K is finite-rank and \mathcal{H} is finite-dimensional. A notable case is that of the RKHSs induced by truncated kernels, i.e. such that there exists d such that $K_{ii} = 0 \forall i > d$. This kind of kernels induce finite-dimensional RKHSs containing FIR systems of order d .*

4.2 Connection with impulse response estimation using orthonormal bases of ℓ_2

As mentioned in Introduction, important impulse response models exploit orthonormal functions $\{\rho_i\}$ in ℓ_2 given e.g. by Laguerre or Kautz models [37]. Then, linear least squares estimators are often adopted to recover the expansion coefficients $\{a_i\}$. Specifically, let $L_k[f]$ be the system output, i.e. the convolution between the known input and f , at the instant t_k where the noisy measurement y_k is available. Then, the impulse response estimate from a data set of size N is

$$\hat{f} = \sum_{i=1}^d \hat{a}_i \rho_i \quad (6a)$$

$$\{\hat{a}_i\}_{i=1}^d = \arg \min_{\{a_i\}_{i=1}^d} \sum_{k=1}^N \left(y_k - L_k \left[\sum_{i=1}^d a_i \rho_i \right] \right)^2 \quad (6b)$$

where d determines model complexity and is typically selected using AIC or cross validation (CV) [34].

An alternative option originally proposed in [45] consists of searching for the impulse response estimate in a stable and infinite-dimensional RKHS with ill-posedness faced by regularization. The least squares estimator (6) is replaced by the following regularized least squares (ReLS) problem

$$\hat{f} = \arg \min_{f \in \mathcal{H}} \sum_{k=1}^N (y_k - L_k[f])^2 + \gamma \|f\|_{\mathcal{H}}^2 \quad (7)$$

where $\|\cdot\|_{\mathcal{H}}$ is the RKHS norm and the positive scalar γ is the so-called regularization parameter. It can e.g. be estimated using empirical Bayes approaches, e.g. see [43].

The results obtained in the previous subsection permit to understand analogies and differences between (6) and (7). In fact, by using Proposition 6 and recalling from [21] also that

$$f = \sum_{i=1}^{\infty} a_i \rho_i \implies \|f\|_{\mathcal{H}}^2 = \sum_{i=1}^{\infty} \frac{a_i^2}{\lambda_i},$$

the following result holds.

Proposition 8 (Representation of ReLS in stable RKHSs) *Let K be stable. Assume that any kernel eigenvalue satisfies $\lambda_i >$*

0 and consider the representation (5) of the induced RKHS. Then, (7) is equivalent to

$$\hat{f} = \sum_{i=1}^{\infty} \hat{a}_i \rho_i \quad (8a)$$

$$\{\hat{a}_i\}_{i=1}^{\infty} = \arg \min_{\{a_i\}_{i=1}^{\infty}} \sum_{k=1}^N \left(y_k - L_k \left[\sum_{i=1}^{\infty} a_i \rho_i \right] \right)^2 + \gamma \sum_{i=1}^{\infty} \frac{a_i^2}{\lambda_i}. \quad (8b)$$

■

So, regularized least squares in a stable (infinite-dimensional) RKHS always model impulse responses using an ℓ_2 orthonormal basis, as in the classical works [55,30]. But the key difference between (6) and (8) is that complexity is no more controlled by the model order since d is set to ∞ . It instead depends on the regularization parameter γ that trades-off data fit and the penalty term. This latter induces stability by constraining the decay rate of the expansion coefficients to zero through the kernel eigenvalues λ_i .

The estimator (7) would seem a computationally unfeasible (infinite-dimensional) variational problem. Actually, according to the representer theorem [31,51,1], the estimate \hat{f} belongs to a subspace of dimension equal to the data-set size N . For dynamic systems, it is determined by the kernel and the system input, see [46][Part III] for details. The kernel implicit encoding so permits to compute the impulse response estimate without knowing the basis functions $\{\rho_i\}$. However, achieving the N expansion coefficients requires $O(N^3)$ operations, so that for large data sets alternative procedures are desirable. A strategy for approximating (7) is to use the equivalence with (8) then resorting to truncated Mercer expansions. Specifically, Problem (8) is replaced by the following d -dimensional surrogate

$$\hat{f}^{(d)} = \sum_{i=1}^d \hat{a}_i^{(d)} \rho_i \quad (9a)$$

$$\{\hat{a}_i^{(d)}\}_{i=1}^d = \arg \min_{\{a_i\}_{i=1}^d} \sum_{k=1}^N \left(y_k - L_k \left[\sum_{i=1}^d a_i \rho_i \right] \right)^2 + \gamma \sum_{i=1}^d \frac{a_i^2}{\lambda_i}. \quad (9b)$$

Note that d has not to trade-off bias and variance here as it instead happens in (6). It has instead to be sufficiently large so that $\hat{f}^{(d)}$ is close to \hat{f} . Indeed, in [41] it has been shown that convergence holds in the RKHS norm as d grows to infinity. In addition, in [41,10] numerical experiments have shown that a relatively small number of eigenfunctions (w.r.t. the data set size N) can provide really good approximations. This is advantageous since, after numerically computing each value $L_k[\rho_i]$, the estimate $\hat{f}^{(d)}$ requires only $O(Nd^2)$ operations.

4.3 Numerical recovery of the Mercer ℓ_2 feature map

In the previous subsection, we have outlined that Mercer expansions of K can be important also for implementing ReLS. However, obtaining closed form expressions of the Mercer feature map is often prohibitive. The following result fills this gap by showing that the ℓ_2 basis of a stable RKHS and the kernel eigenvalues can be numerically estimated (with arbitrary precision) by a sequence of SVDs applied to truncated kernels. This result is not trivial since, in the literature, the problem could not even be posed on a firm theoretical ground. In fact, it was not known whether a stable kernel admitted a Mercer expansion in ℓ_2 , a fact now established by Proposition 5.

Given a kernel K , the notation $K^{(d)}$ indicates its truncated version, i.e. the $d \times d$ matrix obtained by retaining only its first d rows and columns. Then, $\rho_i^{(d)}$ and $\lambda_i^{(d)}$ are the eigenvectors (thought of as elements of ℓ_2 with a tail of zeros) and the eigenvalues obtained by the SVD of $K^{(d)}$. Single multiplicity is assumed for each λ_i , see Remark 21 in Appendix for further discussions.

Theorem 9 (Estimation of Mercer expansions in ℓ_2) *Let K be stable or, more generally, be a kernel inducing a compact operator as a map from ℓ_2 into ℓ_2 itself. Let also ρ_i and λ_i denote, respectively, its eigenfunctions (forming an orthonormal basis in ℓ_2) and the corresponding eigenvalues. Assume also that the multiplicity of each λ_i is equal to one. Then, for any i , as*

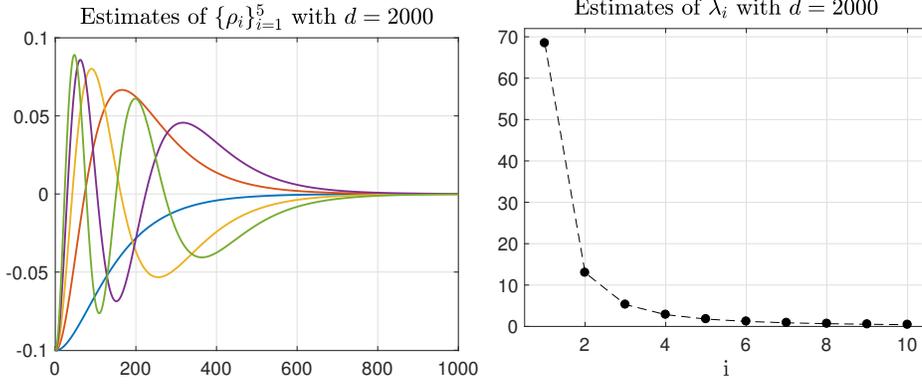


Fig. 2. Estimates of the first 5 eigenfunctions (left) and of the first 10 eigenvalues (right) achieved by the SVD of the truncated stable spline kernel with $d = 2000$.

d grows to ∞ it holds that

$$\lambda_i^{(d)} \rightarrow \lambda_i \tag{10a}$$

$$\|\rho_i^{(d)} - \rho_i\|_2 \rightarrow 0. \tag{10b}$$

where $\|\cdot\|_2$ is the ℓ_2 -norm. ■

Theorem 9 is now applied to the first-order stable spline (SS) kernel [45], also called TC kernel in [17]. This model is often used to describe smooth and exponentially decaying impulse responses. Its (i, j) entry is

$$K_{ij} = \alpha^{\max(i,j)} \tag{11}$$

where the scalar $0 \leq \alpha < 1$ regulates the decay rate of the functions contained in the induced RKHS. In what follows, we set $\alpha = 0.95$.

Our aim is to obtain a good approximation of the SS Mercer expansion in ℓ_2 . For this purpose, we exploit Theorem 9 by computing SVDs of truncated SS kernels of size $d = 200, 400, \dots, 2000$. Both λ_i and $\lambda_i^{(d)}$ are ordered in decreasing order in what follows.

Fig. 2 plots the estimates of the first 5 eigenfunctions (left) and of the first 10 eigenvalues (right) achieved with $d = 2000$. The capability of the finite-dimensional estimator (9) to approximate (7) for small values of d will depend on the number of data available, the system input and the value of the regularization parameter γ . However, the fact that the eigenvalues profile shows that most of the energy of the TC kernel is captured by the first 5-10 eigenfunctions suggests that values of d much smaller than N can do a good job, confirming the experimental results reported in [10].

Fig. 3 provides some details on the reconstruction of the 100-th eigenfunction as d increases from 200 to 2000. The left panel shows the following ℓ_2 norms

$$\left\| \rho_{100}^{(200k+200)} - \rho_{100}^{(200k)} \right\|_2,$$

as a function of the integer k . Such norms can be monitored to assess the convergence (ensured by Theorem 9) of the $\rho_{100}^{(d)}$ towards the eigenfunction ρ_{100} of K . One can see that for values of $k > 4$ the discrepancy quickly goes to zero. The right panel finally plots the approximation of ρ_{100} provided by $\rho_{100}^{(2000)}$.

5 The necessary and sufficient condition for RKHS stability in the Mercer feature space

So far, our starting point has been a kernel designed by specifying its entries K_{ij} . This modeling approach translates the expected features of an impulse response into kernel properties, e.g. smooth exponential decay as described by (11). This way takes advantage of basis functions implicit encoding. Recent literature has shown that also models built by designing eigenfunctions ρ_i and eigenvalues λ_i are valuable. In fact, kernels relying on Laguerre or Kautz functions, that belong to the more general class of Takenaka-Malmquist orthogonal basis functions [30], are useful to describe oscillatory behavior or

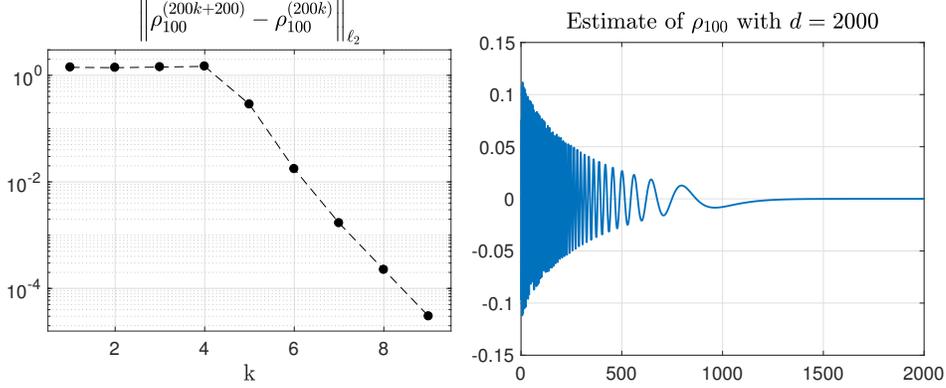


Fig. 3. Distances measured by the ℓ_2 norm between $\rho_{100}^{(200k+200)}$ and $\rho_{100}^{(200k)}$ for different k values (left) and estimate of ρ_{100} obtained with $d = 2000$ (right).

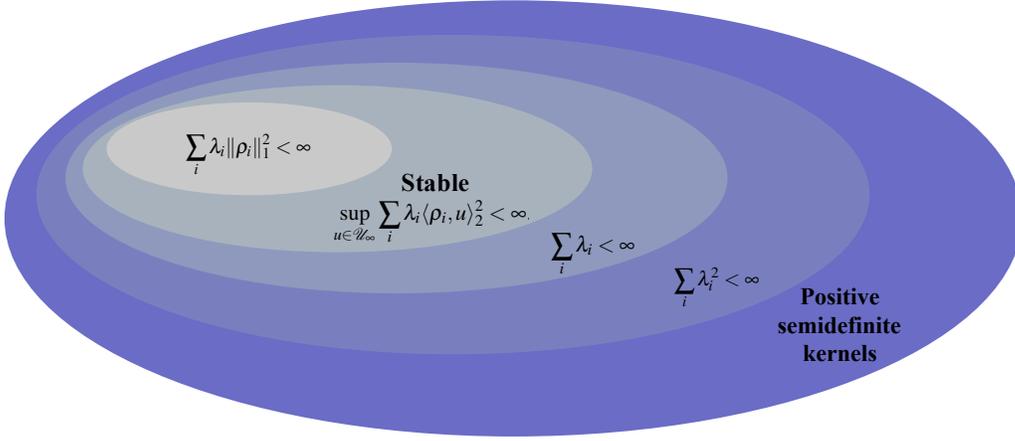


Fig. 4. Inclusion properties of different kernel classes in terms of the Mercer feature space. This representation is the dual of that reported in Fig. 1. Here the kernel sets are defined through properties of the kernel eigenvectors ρ_i , forming an orthonormal basis in ℓ_2 , and of the corresponding kernel eigenvalues λ_i . The condition $\sum_i \lambda_i \|\rho_i\|_1^2 < \infty$ is the most restrictive since it implies kernel absolute summability. The necessary and sufficient condition for stability is $\sup_{u \in \mathcal{U}_\infty} \sum_i \lambda_i \langle \rho_i, u \rangle_2^2 < \infty$. Finally, $\sum_i \lambda_i < \infty$ and $\sum_i \lambda_i^2 < \infty$ are exactly the conditions for a kernel to be finite-trace and squared summable, respectively, see also the proof of Theorem 3 in Appendix for details.

presence of fast/slow poles.

This fact motivates the following fundamental problem. Assigned an orthonormal basis $\{\rho_i\}$ of ℓ_2 , e.g. of the Takenaka-Malmquist type, which conditions on the eigenvalues λ_i ensure that the kernel $K_{xy} = \sum_{i=1}^{+\infty} \lambda_i \rho_i(x) \rho_i(y)$ is stable? One should also expect that, if $\lambda_i > 0 \forall i$, there exist bases that never satisfy such requirement. All of these issues find a definite answer in the following result that provides the necessary and sufficient condition for kernel stability starting from Mercer expansions.

Theorem 10 (RKHS stability using Mercer feature maps) *Let \mathcal{H} be the RKHS induced by K having Mercer expansion $K_{xy} = \sum_{i=1}^{+\infty} \lambda_i \rho_i(x) \rho_i(y)$ with $\{\rho_i\}$ an orthonormal basis of ℓ_2 . Define also*

$$\mathcal{U}_\infty = \left\{ u \in \ell_\infty : |u(i)| = 1, \forall i \geq 1 \right\}.$$

Then, it holds that

$$\mathcal{H} \subset \ell_1 \iff \sup_{u \in \mathcal{U}_\infty} \sum_i \lambda_i \langle \rho_i, u \rangle_2^2 < +\infty \quad (12)$$

where $\langle \cdot, \cdot \rangle_2$ is the inner product in ℓ_2 .

■

We discuss some consequences of the above result.

If there is one ρ_i corresponding to $\lambda_i > 0$ that doesn't belong to ℓ_1 then stability is prevented. In fact $\langle \rho_i, u \rangle_2 = +\infty$ for u containing the signs of the components of ρ_i . Nothing is however required for the eigenvectors associated to $\lambda_i = 0$. Another outcome is the following sufficient stability condition.

Theorem 11 (Sufficient stability condition using Mercer) *Let \mathcal{H} be the RKHS induced by K having Mercer expansion $K_{xy} = \sum_{i=1}^{+\infty} \lambda_i \rho_i(x) \rho_i(y)$ with $\{\rho_i\}$ an orthonormal basis of ℓ_2 . One has*

$$\mathcal{H} \subset \ell_1 \iff \sum_i \lambda_i \|\rho_i\|_1^2 < +\infty. \quad (13)$$

Moreover, such condition also implies kernel absolute summability and, hence, it is not necessary for RKHS stability. ■

The stability condition (13) can be easily used to design model stable impulse responses starting from any kind of basis in ℓ_2 . For instance, as we have recalled in Introduction, Laguerre or Kautz basis functions are often used to define $\{\rho_i\}$ that embed information on the dominant pole of the system and/or presence of resonances. Once such a basis is fixed, one thus assumes that the impulse response has the form

$$f = \sum_{i=1}^{\infty} a_i \rho_i.$$

Now, to exploit the regularized estimator (8) to identify the system, the key point is to understand which kind of constraints on the a_i lead to stable models. The question is equivalent to understand which decay rate of the λ_i ensures that the regularizer

$$\sum_{i=1}^{\infty} \frac{a_i^2}{\lambda_i}$$

enforces impulse responses' absolute summability. As a concrete example, Laguerre and Kautz models all belong to the more general Takenaka-Malmquist class of basis functions known to satisfy the constraint

$$\|\rho_i\|_1 \leq A i,$$

with A a constant independent of i , e.g. see [30]. Theorem 11 then allows us to immediately conclude that the choice

$$\lambda_i \propto i^\nu, \quad \nu > 2$$

always enforces stability in the estimation process for all the Takenaka-Malmquist class. So, any stable impulse response model relying e.g. on Laguerre or Kautz can now embed such a constraint on the eigenvalues' decay.

If the orthonormal basis functions corresponding to strictly positive eigenvalues are all contained in a ball of ℓ_1 , the following result holds.

Theorem 12 (Stability with bases uniformly bounded in ℓ_1) *Let K be a kernel having Mercer expansion $K_{xy} = \sum_{i=1}^{+\infty} \lambda_i \rho_i(x) \rho_i(y)$ with $\{\rho_i\}$ an orthonormal basis of ℓ_2 and $\|\rho_i\|_1 \leq A < +\infty$ if $\lambda_i > 0$ (A is a constant independent of i). Then, one has*

$$\mathcal{H} \subset \ell_1 \iff \sum_i \lambda_i < +\infty. \quad (14)$$

■

Finally, all the new insights on kernel stability in the Mercer feature space are graphically depicted in Fig. 4.

6 Conclusions

The results reported in this paper shed new light on the RKHSs containing absolutely summable impulse responses. The inclusion properties here derived give a clear picture on the relationship between stable kernels and other fundamental classes. They

have important consequences for stability tests. In addition, they provide representations of RKHSs and of related regularized least squares estimators that clarify the relationship with linear system identification via orthonormal bases in ℓ_2 . Our analysis includes also the necessary and sufficient stability condition for kernels built through these functions. The paper thus provides new mathematical foundations of stable RKHSs with impact also on stable impulse responses modeling and estimation.

7 Appendix

In what follows, given a finite- or infinite-dimensional matrix A , the notation $A \succeq 0$ will indicate that the matrix is symmetric and positive semidefinite. Moreover, the canonical basis in ℓ_2 will be denoted by $\{e_i\}, i \in \mathbb{N}$.

We will also often use M_m to denote a matrix of size $m \times m$. In addition, let

$$\|M_m\|_{\infty,1} := \max_{\|u\|_{\infty}=1} \|M_m u\|_1, \quad (15)$$

that corresponds to the norm of the linear operator $M_m : \mathbb{R}^m \rightarrow \mathbb{R}^m$ with the domain and co-domain equipped, respectively, with the ℓ_{∞} and the ℓ_1 norms.

In this Appendix, we will also adopt a different notation for a kernel and a kernel operator (so far indicated indistinctly with K). The notation M will indicate an infinite-dimensional matrix representing a kernel, i.e. $M \succeq 0$. The associated kernel operator is \mathcal{M} . This will thus be a self-adjoint positive semidefinite operator with domain and co-domain specified later on.

In addition, given any integer $r \geq 1$, with possibly also $r = \infty$, the set \mathcal{U}_r is defined as follows

$$\mathcal{U}_r := \{x \in \mathbb{R}^r : x(i) = \pm 1, \forall i = 1, \dots, r\}. \quad (16)$$

7.1 Proof of Theorem 3

Let p an integer ($p \geq 1$) that also defines the odd number $m = 2p + 1$ and the corresponding power of two $n = 2^m$. Let also $x_i \in \mathcal{U}_m$ ($i = 1, 2, \dots, n$) that, according to (16), are distinct vectors containing exactly m elements ± 1 (the ordering of such vectors is irrelevant). Then, for any $n = 2^3, 2^5, 2^7, \dots, V^{(n)}$ indicates the matrix of size $n \times m$ given by

$$V^{(n)} = \begin{bmatrix} x_1 & x_2 & \dots & x_n \end{bmatrix}^{\top}. \quad (17)$$

Thus, the rows of such matrices contain all the possible permutations of ± 1 . As an example, setting $p = 1$ and, hence, $m = 3, n = 2^3 = 8$, one obtains

$$V^{(8)} = \begin{pmatrix} 1 & 1 & 1 \\ 1 & 1 & -1 \\ 1 & -1 & 1 \\ 1 & -1 & -1 \\ -1 & 1 & 1 \\ -1 & 1 & -1 \\ -1 & -1 & 1 \\ -1 & -1 & -1 \end{pmatrix} \quad (18)$$

1. $\mathcal{S}_1 \subset \mathcal{S}_s$

The inclusion $\mathcal{S}_1 \subseteq \mathcal{S}_s$ is immediate and well-known in the literature, as also discussed after stating Theorem 3. The (strict) inclusion $\mathcal{S}_1 \subset \mathcal{S}_s$ is not trivial and its proof can be found in [8]. It relies on the building of a particular kernel, depending on the matrices $V^{(n)}$ defined in (17), that is stable but non absolutely summable.

2. $\mathcal{S}_s \subset \mathcal{S}_{ft}$

Lemma 13 Let $M_m = M_m^T \succeq 0$ of size $m \times m$. Then the following inequalities hold true

$$\text{tr}(M_m) \leq \|M_m\|_{\infty,1} \leq n \text{tr}(M_m)$$

Proof: Using the same arguments contained in the proof of Lemma 3.1 in [8], it holds that

$$\|M_m\|_{\infty,1} = \max_{x \in \mathcal{U}_m} \|M_m x\|_1$$

Recalling that $V^{(n)\top}$ contains all the vectors in \mathcal{U}_m as columns, the problem corresponds to evaluating

$$M_m V^{(n)\top}$$

and looking for the column with maximum ℓ_1 norm. The ℓ_1 norm of any column is easily obtained by means of a scalar product of the column itself with a suitable $x \in \mathcal{U}_m$ corresponding to the signs of the column entries. So, one has that

$$V^{(n)} M_m V^{(n)\top}$$

surely contains (within its n^2 entries) these n ℓ_1 norms. Also, the searched maximum ℓ_1 norm coincides with the maximum of all n^2 entries since $x_1^T c \leq x_2^T c$, $\forall x_1 \in \mathcal{U}_m$ if $x_2 = \text{sign}(c)$ (where, for each entry of c , sign returns 1 if it is larger than zero and -1 otherwise). Now, $V^{(n)} M_m V^{(n)\top} \succeq 0$, which implies that its maximum entry appears along its diagonal, so

$$\|M_m\|_{\infty,1} = \max_{i=1,\dots,n} [V^{(n)} M_m V^{(n)\top}]_{ii}$$

Now, note that the trace of $V^{(n)} M_m V^{(n)\top}$ satisfies

$$\text{tr}[V^{(n)} M_m V^{(n)\top}] \geq \|M_m\|_{\infty,1} \geq \frac{1}{n} \text{tr}[V^{(n)} M_m V^{(n)\top}].$$

Finally,

$$\begin{aligned} \text{tr}[V^{(n)} M_m V^{(n)\top}] &= \text{tr}[M_m V^{(n)\top} V^{(n)}] \\ &= \text{tr}[M_m (nI_m)] = n \text{tr}[M_m] \end{aligned}$$

and this concludes the proof. □

Lemma 14 Let $\mathcal{M} : \ell_\infty \rightarrow \ell_1$ be a self-adjoint, positive semidefinite and bounded operator. Then \mathcal{M} is finite-trace, too.

Proof: By denoting with M_k the $k \times k$ submatrix of the infinite matrix that represents the operator \mathcal{M} , it is easy to see that

$$\|M_k\|_{\infty,1} \leq \|\cdot\mathcal{M}\|_{\infty,1} < +\infty, \quad \forall k = 1, 2, \dots,$$

where $\|\cdot\mathcal{M}\|_{\infty,1}$ is the operator norm of \mathcal{M} . So, exploiting Lemma 13, we also have

$$\text{tr}[M_k] \leq \|\cdot\mathcal{M}\|_{\infty,1}, \quad \forall k = 1, 2, \dots$$

Then, since $\text{tr}[M_k]$ is a monotone non-decreasing sequence upper-bounded by $\|\cdot\mathcal{M}\|_{\infty,1} < +\infty$, letting $M_k(k, k)$ the k -th element along the diagonal of M_k , this implies that

$$\text{tr}[\cdot\mathcal{M}] := \sum_{k=1}^{+\infty} M_k(k, k) \leq \|\cdot\mathcal{M}\|_{\infty,1} < +\infty.$$

□

Lemma 14 thus shows that $\mathcal{S}_s \subseteq \mathcal{S}_{ft}$. To prove that the inclusion is strict, it suffices to consider the infinite-dimensional matrix (kernel)

$$M = vv^T.$$

In fact, one has $\text{tr}(M) = \|v\|_2^2 < +\infty$ iff $v \in \ell_2$. If $v \notin \ell_2$, letting $w = \text{sign}(v) \in \ell_\infty$ one obtains $Mw = v\|v\|_1 = \infty$. This proves that the kernel M is unstable.

3. $\mathcal{S}_{ft} \subset \mathcal{S}_2$

Given a kernel, represented by an infinite-dimensional matrix $M = M^T \succeq 0$, it is now important to consider the induced kernel operator \mathcal{M} as a map from ℓ_2 into itself. Its operator norm is given by

$$\|\mathcal{M}\|_2 = \sup_{\|v\|_2=1} \|\mathcal{M}v\|_2.$$

In addition, given any orthonormal basis $\{v_i\}$, its nuclear norm is

$$\sum_{i \geq 1} \langle v_i, \mathcal{M}v_i \rangle, \quad (19)$$

while its (squared) Hilbert-Schmidt (HS) norm is

$$\sum_{i=1}^{\infty} \|\mathcal{M}v_i\|_2^2 \quad (20)$$

and both of them turn out independent of the particular chosen basis.

This view allows us to cast the relation between \mathcal{S}_{ft} and \mathcal{S}_2 in terms of the important nuclear and HS operators. First, we briefly recall some fundamental results that can be found e.g. in [33,26,49]. By definition, an operator is HS if (20) is finite. In particular, a kernel operator \mathcal{M} is HS iff it is induced by a squared summable kernel. An HS kernel operator is always compact with squared norm (20) given by

$$\sum_{i \geq 1} \lambda_i^2 \quad (21)$$

where the λ_i are the eigenvalues of M .

An operator is said to be nuclear if it can be written as the composition of two HS operators. Any nuclear operator is HS. Also, a positive semidefinite operator is nuclear if and only if it is compact with

$$\sum_{i \geq 1} \lambda_i < +\infty, \quad (22)$$

where the λ_i still denote the eigenvalues of M [49][Section 2].

Now, we want to prove that all the operators induced by the kernels in \mathcal{S}_{ft} are nuclear. For any $n \geq 0$, define the projection operators \mathcal{P}_n and the partial traces $\text{tr}_n(\mathcal{M})$ as follows

$$\mathcal{P}_n : \sum_{h=1}^{+\infty} a_h e_h \rightarrow \sum_{h=1}^n a_h e_h$$

$$\text{tr}_n(\mathcal{M}) := \sum_{h=n+1}^{+\infty} M(h, h).$$

Lemma 15 *Let $\mathcal{P}_n v = 0$. Then*

$$\|\mathcal{M}v\|_2 \leq \sqrt{\text{tr}_n(\mathcal{M})} \sqrt{\text{tr}(\mathcal{M})} \|v\|_2$$

Proof: We have $v = \sum_{h=n+1}^{+\infty} a_h e_h$, so

$$\begin{aligned}
\|\mathcal{M}v\|_2^2 &= \sum_{r=1}^{+\infty} \left| \sum_{s=n+1}^{+\infty} M(r,s) a_s \right|^2 \\
&\leq \sum_{r=1}^{+\infty} \left(\sum_{s=n+1}^{+\infty} |M(r,s)| \cdot |a_s| \right)^2 \leq \sum_{r=1}^{+\infty} \left(\sum_{s=n+1}^{+\infty} a_s^2 \cdot \sum_{s=n+1}^{+\infty} M^2(r,s) \right) \\
&= \sum_{s=n+1}^{+\infty} a_s^2 \left(\sum_{r=1}^{+\infty} \sum_{s=n+1}^{+\infty} M^2(r,s) \right) = \|v\|_2^2 \left(\sum_{r=1}^{+\infty} \sum_{s=n+1}^{+\infty} M^2(r,s) \right) \\
&\leq \|v\|_2^2 \left(\sum_{r=1}^{+\infty} \sum_{s=n+1}^{+\infty} M(r,r) M(s,s) \right) = \|v\|_2^2 \sum_{r=1}^{+\infty} W(r,r) \sum_{s=n+1}^{+\infty} W(s,s) \\
&= \text{tr}_n(\mathcal{W}) \text{tr}(\mathcal{W}) \|v\|_2^2
\end{aligned}$$

and this completes the proof. \square

Lemma 16 *Any finite-trace kernel operator from ℓ_2 into itself is compact.*

Proof: Let $(w^{(h)}, h \in \mathbb{N})$ be a bounded sequence of ℓ_2 elements, i.e. $\|w^{(h)}\| \leq A$ for some $A > 0$ and for any $h \in \mathbb{N}$. Consider the sequence $\mathcal{P}_n w^{(h)}$ that belongs to a subspace isomorphic to \mathbb{R}^n . So, a subsequence $(w^{(h)}, h \in \mathcal{I}_n)$ exists such that $(\mathcal{P}_n w^{(h)}, h \in \mathcal{I}_n)$ converges to some $v_n = \sum_{h=1}^n a_h e_h$, with $\|v_n\|_2 \leq A$ by $\|\mathcal{P}_n w^{(h)}\|_2 \leq \|w^{(h)}\|_2 \leq A$. Now we can proceed inductively as follows. The $(\mathcal{P}_{n+1} w^{(h)}, h \in \mathcal{I}_n)$ are equal to $(\mathcal{P}_n w^{(h)}, h \in \mathcal{I}_n)$, except for the $(n+1)$ -entry, which is upper bounded (in absolute value) by A . So, a subsequence $(w^{(h)}, h \in \mathcal{I}_{n+1})$ of the $(w^{(h)}, h \in \mathcal{I}_n)$ exists such that $(\mathcal{P}_{n+1} w^{(h)}, h \in \mathcal{I}_{n+1})$ converges to some $v_{n+1} = \sum_{h=1}^{n+1} a_h e_h$. Note that the first n entries are exactly the same for both v_{n+1} and v_n , so another a_i 's coefficient has been added without modifying the first n coefficients. In this way, we can finally obtain a vector $v \in \ell_2$

$$v = \sum_{h=1}^{+\infty} a_h e_h, \quad \|v\|_2 \leq A$$

which is the limit in ℓ_2 of the sequence v_n , with $\|v_n\|_2$ forming a monotone non-decreasing sequence upper bounded by A . Hence, also $\|v\|_2 \leq A$. We have now

$$\begin{aligned}
\|\mathcal{M}w^{(h)} - \mathcal{M}v\|_2 &= \|\mathcal{M}(\mathcal{I} - \mathcal{P}_n)w^{(h)} \\
&\quad + [\mathcal{M}\mathcal{P}_n w^{(h)} - \mathcal{M}v_n] + \mathcal{M}(v_n - v)\|_2 \\
&\leq \|\mathcal{M}(\mathcal{I} - \mathcal{P}_n)w^{(h)}\|_2 \\
&\quad + \|\mathcal{M}(\mathcal{P}_n w^{(h)} - v_n)\|_2 + \|\mathcal{M}(v_n - v)\|_2
\end{aligned}$$

where $\mathcal{P}_n(\mathcal{I} - \mathcal{P}_n)w^{(h)} = 0$, for $h \in \mathcal{I}_n$, so Lemma 3 applies leading to $\|\mathcal{M}(\mathcal{I} - \mathcal{P}_n)w^{(h)}\|_2 \leq \sqrt{\text{tr}_n(\mathcal{M})} \sqrt{\text{tr}(\mathcal{M})}$, for $h \in \mathcal{I}_n$. Therefore the first and the third term are both infinitesimal w.r.t. n for any $h \in \mathcal{I}_n$. The second term is also infinitesimal w.r.t. $h \in \mathcal{I}_n$ for any n since, using Lemma 15 with $n = 0$, for any $v \in \ell_2$ one has

$$\|\mathcal{M}v\|_2 \leq \text{tr}(\mathcal{M}) \|v\|_2.$$

Let now ε_k be any monotone non-increasing sequence converging to zero, and let $n(k)$ be such that the first and the third term are both less than $\frac{\varepsilon_k}{3}$ for any $h \in \mathcal{I}_{n(k)}$. Let also $h(k) \in \mathcal{I}_{n(k)}$ be such that the second term is less than $\frac{\varepsilon_k}{3}$. Thus, we have

$$0 \leq \|\mathcal{M}w^{(h(k))} - \mathcal{M}v\|_2 < \varepsilon_k.$$

In this inductive procedure w.r.t. k , we only need to choose $h(1) < h(2) < \dots < h(k) < \dots$ and this is always possible since $h(k)$ can be chosen in infinitely many ways in view of the convergence property of the second term w.r.t. h . Finally, by defining the countable set $\mathcal{I} := \{h(1), h(2), \dots, h(k), \dots\}$, the subsequence $(w^{(h)}, h \in \mathcal{I})$ of the original sequence $(w^{(h)}, h \in \mathbb{N})$, satisfies

$$0 \leq \|\mathcal{M}w^{(h(k))} - \mathcal{M}v\|_2 < \varepsilon_k, \quad \forall k = 1, 2, \dots$$

with $h(k)$ strictly monotone increasing. Since ε_k is infinitesimal, $\mathcal{M}w^{(h(k))}$ converges to $\mathcal{M}v$. So, any bounded sequence $w^{(h)} \in \ell_2$ admits a subsequence $w^{h(k)}$ such that $\mathcal{M}w^{h(k)}$ is convergent in ℓ_2 , proving the compactness of the operator [50,58]. \square

Combination of the Lemma 16 and of the spectral theorem [26] ensures that there exists a complete orthonormal basis of M given by eigenvectors $\{\rho_i\}$ of M with corresponding eigenvalues denoted by $\{\lambda_i\}$. Using first the $\{\rho_i\}$ and then the canonical basis $\{e_i\}$ of ℓ_2 to evaluate the nuclear norm (19), if $M \in \mathcal{S}_{ft}$ one obtains

$$\sum_{i \geq 1} \lambda_i = \sum_{i \geq 1} M_{ii} < +\infty.$$

So, (22) holds true and \mathcal{M} is a nuclear operator. Then, the set inclusion immediately derives from the fact that any nuclear operator is also HS. Such inclusion is obviously strict as the simple example $M = \text{diag}\{1, 1/2, 1/3, \dots, 1/k, \dots\}$ shows.

4. $\mathcal{S}_2 \subset \text{Kernels set}$

The kernels set contains all the positive semidefinite infinite matrices. The inclusion is then obvious and is strict as proved by the example vv^T with all the entries of the infinite-dimensional column vector v equal to 1.

7.2 Proof of Theorem 9

Given the infinite-dimensional matrix $M = M^T \succeq 0$ associated with a compact operator, let M_d contain the first d rows and columns of M . Then, we will consider the following partition

$$M = \begin{bmatrix} M_d & A_d \\ A_d^T & B_d \end{bmatrix}$$

with the eigenvalues of M and M_d denoted, respectively, by

$$\lambda_1(M) \geq \lambda_2(M) \geq \dots \quad \text{and} \quad \lambda_1(d) \geq \lambda_2(d) \geq \dots$$

For the moment, no assumption on eigenvalues multiplicities is used. In addition, $\langle \cdot, \cdot \rangle_2$ denotes the inner-product in ℓ_2 for infinite-dimensional vectors or in the classical Euclidean space for finite-dimensional ones. The same holds for $\|\cdot\|_2$.

Lemma 17 For any $d \geq k \geq 1$ it holds that

$$\max_{\langle v_h, v_k \rangle = \delta_{hk}} \sum_{h=1}^k v_h^T M v_h = \sum_{h=1}^k \lambda_h(M), \tag{23a}$$

$$\max_{\langle v_h, v_k \rangle = \delta_{hk}} \sum_{h=1}^k v_h^T M_d v_h = \sum_{h=1}^k \lambda_h(d) \tag{23b}$$

Proof: We will exploit the spectral theorem that holds true both for M and for M_d . Using an orthonormal basis, either in ℓ_2 or in \mathbb{R}^d , the two equalities are obtained by choosing v_h as (one of) the eigenvectors corresponding to either $\lambda_h(M)$ or $\lambda_h(d)$. Now, it suffices to prove that any other choice of the v_h does not lead to results larger than the sum of the first k eigenvalues. We can just focus on M . Assume that $v_h = \sum_{i=1}^{+\infty} a_{hi} \rho_i$, with $h = 1, 2, \dots, k$, are orthonormal vectors, expressed in terms of the orthonormal basis ρ_i (each ρ_i is associated with $\lambda_i(M)$). Let also v_{h+1}, v_{h+2}, \dots be any completion of the set $\{v_h, h = 1, 2, \dots, k\}$ up to an orthonormal ℓ_2 basis. Since a_{hi} , with $h = 1, 2, \dots, k$, are the first k elements of the i -th row of a unitary (infinite) matrix U with columns given by the vectors v_h , one has $\sum_{h=1}^k a_{hi}^2 \leq 1 \forall i \in \mathbb{N}$ so that

$$\begin{aligned} \sum_{h=1}^k v_h^T M v_h &= \sum_{i=1}^{+\infty} \lambda_i \left(\sum_{h=1}^k a_{hi}^2 \right) \leq \sum_{i=1}^{k-1} \lambda_i \left(\sum_{h=1}^k a_{hi}^2 \right) + \lambda_k \left(\sum_{i=k}^{+\infty} \sum_{h=1}^k a_{hi}^2 \right) \\ &= \sum_{i=1}^{k-1} \lambda_i \left(\sum_{h=1}^k a_{hi}^2 \right) + \lambda_k \left(\sum_{h=1}^k \sum_{i=k}^{+\infty} a_{hi}^2 \right) = \sum_{i=1}^{k-1} \lambda_i \left(\sum_{h=1}^k a_{hi}^2 \right) + k\lambda_k - \sum_{i=1}^{k-1} \lambda_k \left(\sum_{h=1}^k a_{hi}^2 \right) \\ &= \sum_{i=1}^{k-1} (\lambda_i - \lambda_k) \left(\sum_{h=1}^k a_{hi}^2 \right) + k\lambda_k \leq \sum_{i=1}^{k-1} (\lambda_i - \lambda_k) + k\lambda_k = \sum_{i=1}^k \lambda_i \end{aligned}$$

which completes the proof. \square

Lemma 18 *One has*

$$\|A_d w\|_2^2 \leq \lambda_1(d) \lambda_M(B_d) \|w\|_2^2 \quad (24)$$

where $\lambda_M(B_d)$ is the maximum eigenvalue of B_d .

Proof: From

$$\begin{bmatrix} v^T & w^T \end{bmatrix} \begin{bmatrix} M_d & A_d \\ A_d^T & B_d \end{bmatrix} \begin{bmatrix} v \\ w \end{bmatrix} \geq 0$$

it follows that $|v^T A_d w|^2 \leq (v^T M_d v)(w^T B_d w)$. Now, by choosing $v = \frac{A_d w}{\|A_d w\|_2}$, the previous inequality becomes $\|A_d w\|_2^2 \leq \lambda_1(d) \lambda_M(B_d) \|w\|_2^2$ where we have applied (23) with $k = 1$ to M_d and B_d . \square

Lemma 19 *One has*

$$\lim_{d \rightarrow \infty} \lambda_M(B_d) = 0 \quad (25)$$

Proof: Let z_d be a unit norm eigenvector of B_d corresponding to $\lambda_M(B_d)$, and define $q_d = \begin{bmatrix} 0_d & z_d^T \end{bmatrix}^T$ where 0_d is a column vector with d zero entries. We have $\|q_d\|_2 = 1$ by construction, and $M q_d = \begin{bmatrix} A_d z_d \\ B_d z_d \end{bmatrix} = \begin{bmatrix} A_d z_d \\ \lambda_M(B_d) z_d \end{bmatrix}$ with $\|M q_d\|_2 \geq \lambda_M(B_d)$. By compactness, there exists $d(k)$ s.t. $M q_{d(k)} \rightarrow q \in \ell_2$, so that $\|\lambda_M(B_{d(k)}) q_{d(k)} - (\mathcal{I} - \mathcal{P}_{d(k)}) q\|_2$ is converging to zero. In terms of squared norms, this implies that $\lambda_M^2(B_{d(k)}) - \sum_{h=d(k)+1}^{+\infty} q^2(h)$ goes to zero but $\sum_{h=d(k)+1}^{+\infty} q^2(h)$ is also infinitesimal proving that $\lambda_M(B_{d(k)}) \rightarrow 0$. Since $\lambda_M(B_d)$ is a monotone non-increasing sequence (a fact that can be proved using (23) with $k = 1$ and M replaced by B_d), $\lambda_M(B_d)$ is infinitesimal too. \square

Lemma 20 *For any $k \geq 1$, one has*

$$\lim_{d \rightarrow +\infty} \max_{\langle s_h, s_k \rangle_2 = \delta_{hk}} \sum_{h=1}^k s_h^T M_d s_h = \max_{\langle s_h, s_k \rangle_2 = \delta_{hk}} \sum_{h=1}^k s_h^T M s_h. \quad (26)$$

Proof: As shown in the proof of Lemma 17, the maximum on the r.h.s. exists. Let it be attained for some (orthonormal) vectors $\begin{bmatrix} v_h^T & z_h^T \end{bmatrix}^T$ where v_h has dimension d and $h = 1, 2, \dots, k$. One thus has

$$\begin{aligned} \max_{\langle s_h, s_k \rangle_2 = \delta_{hk}} \sum_{h=1}^k s_h^T M s_h &= \sum_{h=1}^k \begin{bmatrix} v_h^T & z_h^T \end{bmatrix} \begin{bmatrix} M_d & A_d \\ A_d^T & B_d \end{bmatrix} \begin{bmatrix} v_h \\ z_h \end{bmatrix} \\ &= \sum_{h=1}^k v_h^T M_d v_h + 2 \sum_{h=1}^k v_h^T A_d z_h + \sum_{h=1}^k z_h^T B_d z_h. \end{aligned}$$

As d grows to ∞ , $\|z_h\|_2$ tends to zero while $\|v_h\|_2$ tends to 1. In addition, since $z_h^T B_d z_h \leq \lambda_M(B_d) \|z_h\|_2^2$, from (24), (25) and the inequality $z_h^T B_d z_h \leq \lambda_M(B_d) \|z_h\|_2^2$ it comes out that $v_h^T A_d z_h$ and $z_h^T B_d z_h$ also go to zero. From the orthonormality constraint $\langle v_i, v_j \rangle_2 = \delta_{ij} - \langle u_i, u_j \rangle_2$, one then has that the vectors v_i tend to become mutually orthonormal. In particular, by applying the Gram-Schmidt orthonormalization procedure to the v_i , one obtains $v_i = w_i + \varepsilon_i$, with the w_i mutually orthonormal and the $\|\varepsilon_i\|_2$

tending to zero. So, the term $v_h^T M_d v_h$ can be written as $w_h^T M_d w_h + \delta_h$, with δ_h converging to zero. Overall, one has

$$\begin{aligned} \max_{\langle s_h, s_k \rangle = \delta_{hk}} \sum_{h=1}^k s_h^T M s_h &= \lim_{d \rightarrow +\infty} \sum_{h=1}^k w_h^T M_d w_h \\ &\leq \lim_{d \rightarrow +\infty} \max_{\langle s_h, s_k \rangle = \delta_{hk}} \sum_{h=1}^k s_h^T M_d s_h \end{aligned}$$

Now, $\max_{\langle s_h, s_k \rangle = \delta_{hk}} \sum_{h=1}^k s_h^T M_d s_h$ is a monotone non-decreasing sequence since maximizing M_{d+1} w.r.t. vectors with the last entry equal to zero corresponds to maximizing M_d . Hence, this sequence of maximum values is upper bounded by $\max_{\langle s_h, s_k \rangle = \delta_{hk}} \sum_{h=1}^k s_h^T M s_h$. Therefore, we also obtain

$$\max_{\langle s_h, s_k \rangle = \delta_{hk}} \sum_{h=1}^k s_h^T M_d s_h \leq \max_{\langle s_h, s_k \rangle = \delta_{hk}} \sum_{h=1}^k s_h^T M s_h$$

which, together with the previous inequality, shows that

$$\lim_{d \rightarrow +\infty} \max_{\langle s_h, s_k \rangle = \delta_{hk}} \sum_{h=1}^k s_h^T M_d s_h = \max_{\langle s_h, s_k \rangle = \delta_{hk}} \sum_{h=1}^k s_h^T M s_h$$

□

Combining (23) and (26), one obtains

$$\lim_{d \rightarrow +\infty} \sum_{h=1}^k \lambda_h(d) = \sum_{h=1}^k \lambda_h(M)$$

with convergence in a monotone non-decreasing sense. Such relation, evaluated for $k = 1$, implies that $\lambda_1(d)$ tends to $\lambda_1(M)$, and then, for $k = 2$, $\lambda_2(d)$ tends to $\lambda_2(M)$, and so on, inductively.

Let's now consider a unit norm eigenvector corresponding to $\lambda_h(M)$ composed by the subvectors v_d and w_d , i.e.

$$\begin{bmatrix} M_d & A_d \\ A_d^T & B_d \end{bmatrix} \begin{bmatrix} v_d \\ w_d \end{bmatrix} = \lambda_h(M) \begin{bmatrix} v_d \\ w_d \end{bmatrix} \Rightarrow (M_d - \lambda_h(M)I)v_d = -A_d w_d. \quad (27)$$

Let v_d be given in terms of the orthonormal eigenvectors $s_1(d), \dots, s_d(d)$ of M_d associated with $\lambda_1(d) \geq \lambda_2(d) \geq \lambda_3(d) \geq \dots$. So

$$v_d = a_1(d)s_1(d) + \dots + a_d(d)s_d(d)$$

where $\|a(d)\|_2 \leq 1$ since $\|v_d\|_2^2 = 1 - \|w_d\|_2^2 \leq 1$. Plugging such expression of v_d in (27) one obtains

$$\sum_{i=1}^d a_i^2(d) (\lambda_h(M) - \lambda_i(d))^2 = \|A_d w_d\|_2^2$$

with $\|a(d)\|_2 \leq 1$, $\|w_d\|_2 \leq 1$. From (24), (25) and the fact that $\|w_d\|_2$ is converging to zero as d goes to ∞ , one obtains that $\|A_d w_d\|_2^2$ is also infinitesimal w.r.t. d . So one has

$$\sum_{i=1}^d a_i^2(d) (\lambda_h(M) - \lambda_i(d))^2 \rightarrow 0 \quad (28a)$$

$$\sum_{i=1}^d a_i^2(d) \rightarrow 1. \quad (28b)$$

Now, let us assume that the multiplicity of $\lambda_h(M) \neq 0$ is $v_h = 1$. The eigenvalues convergence ensures that we can choose $\varepsilon > 0$ less than an half of the minimum between $\lambda_{h-1}(M) - \lambda_h(M)$ and $\lambda_h(M) - \lambda_{h+1}(M)$ such that, for k fixed and $k \geq h$, there exists N such that $d \geq N$ implies $|\lambda_h(M) - \lambda_h(d)| < \varepsilon$, while $|\lambda_h(M) - \lambda_j(d)| > \varepsilon$ for all $j \neq h$. It follows from (28) that the $a_i(d)$ with $i \neq h$ decay to zero as d goes to ∞ . This implies

$$v_d = a_h(d)s_h(d) + \varepsilon_d, \text{ with } a_h^2(d) \rightarrow 1, \|\varepsilon_d\|_2 \rightarrow 0$$

showing that, as d goes to ∞ , one has $\|\pm s_h(d) - v_d\|_2 \rightarrow 0$ where $\pm s_h(d)$ is the eigenvector (possibly corrected to sign) corresponding to the (only) eigenvalue of M_d which tends to $\lambda_h(M)$.

Remark 21 *If the eigenvalue multiplicity is $v_h > 1$ for some h , the eigenvectors are not well-defined, since there exist infinitely many orthonormal bases for the eigenspace. The v_h -dimensional eigenspace is approximated (in some sense) by the corresponding space generated by the eigenvectors of M_d corresponding to the eigenvalues which tend to the same $\lambda_h(M)$. However, nothing can be said about the behavior of the single $a_i(d)$, since this is strongly related to the choice of the eigenvectors $s_i(d)$. One has thus to consider v_h eigenvectors which tend to lie in a v_h -dimensional eigenspace.*

7.3 Proof of Theorem 10

Let \mathcal{M} be the operator induced by the kernel M , thought of as a map from ℓ_∞ into ℓ_1 . Let $\|\mathcal{M}\|_{\infty,1}$ denote its operator norm. Then, we know from Theorem 1, and subsequent discussions, that the necessary and sufficient condition for the stability of the RKHS induced by M is

$$\|\mathcal{M}\|_{\infty,1} < +\infty. \quad (29)$$

Then, let (λ_i, ρ_i) be the eigenvalues and the eigenvectors orthogonal in ℓ_2 associated with \mathcal{M} . The function

$$f(u) := \|y\|_1 = \sum_i |y(i)| = \sum_i \left| \sum_h M_{ih} u(h) \right|$$

is convex being sums of compositions of absolute values and linear functions. This permits to state that, for any fixed variable $u(h)$, its maximum value is obtained either for $u(h) = +1$ or $u(h) = -1$. By an inductive reasoning we thus obtain

$$\|\mathcal{M}\|_{\infty,1} = \sup_{u \in \mathcal{U}_\infty} f(u) = \sup_{u \in \mathcal{U}_\infty} \sum_i \left| \sum_h M_{ih} u(h) \right|.$$

Now, let $M = UDU^T$, where D is diagonal and contains the eigenvalues of M while the columns of U are the corresponding eigenvectors. Then, we have

$$y = Uw, \quad w = DU^T u$$

and, hence,

$$\begin{aligned} w &= \left[\lambda_1 \langle \rho_1, u \rangle_2 \quad \lambda_2 \langle \rho_2, u \rangle_2 \quad \dots \right]^T \\ y &= \lambda_1 \langle \rho_1, u \rangle_2 \rho_1 + \lambda_2 \langle \rho_2, u \rangle_2 \rho_2 + \dots \end{aligned}$$

To evaluate $\|y\|_1$, we need to consider the scalar product $\langle s(u), y \rangle_2$, where $s(u) = \text{sign}(y)$ (since y depends on u , also $s(u)$ does). In fact, we have

$$h(u) := \|y\|_1 = \sum_h \lambda_h \langle \rho_h, u \rangle_2 \langle \rho_h, s(u) \rangle_2$$

and this implies

$$\begin{aligned} \|\mathcal{M}\|_{\infty,1} &= \sup_{u \in \mathcal{U}_\infty} \sum_h \lambda_h \langle \rho_h, u \rangle_2 \langle \rho_h, s(u) \rangle_2 \\ &= \sup_{u \in \mathcal{U}_\infty} h(u). \end{aligned}$$

Consider also

$$g(u) := \sum_h \lambda_h \langle \rho_h, u \rangle_2^2$$

and define

$$A := \sup_{u \in \mathcal{U}_\infty} \sum_h \lambda_h \langle \rho_h, u \rangle_2^2 = \sup_{u \in \mathcal{U}_\infty} g(u).$$

By definition of $s(u)$, it follows that

$$h(u) \geq g(u) \implies \|\mathcal{M}\|_{\infty,1} \geq A.$$

On the other hand

$$\begin{aligned}
h(u) &= \sum_h \lambda_h \langle \rho_h, u \rangle_2 \langle \rho_h, s(u) \rangle_2 \\
&= \sum_h \left(\sqrt{\lambda_h} \langle \rho_h, u \rangle_2 \right) \left(\sqrt{\lambda_h} \langle \rho_h, s(u) \rangle_2 \right) \\
&\leq \sqrt{\sum_h \lambda_h \langle \rho_h, u \rangle_2^2} \sqrt{\sum_h \lambda_h \langle \rho_h, s(u) \rangle_2^2} \\
&\leq \sqrt{g(u)} \sqrt{g(s(u))}
\end{aligned}$$

that implies

$$\|\mathcal{M}\|_{\infty,1} \leq A.$$

So, one has

$$\|\mathcal{M}\|_{\infty,1} = \sup_{u \in \mathcal{U}_\infty} \sum_h \lambda_h \langle \rho_h, u \rangle_2^2.$$

and this, in view of the necessary and sufficient stability condition (29), concludes the proof.

7.4 Proof of Theorem 11

Let again (λ_i, ρ_i) be the eigenvalues and the eigenvectors orthogonal in ℓ_2 associated with \mathcal{M} , the kernel operator induced by M . One has

$$\begin{aligned}
|\langle \rho_h, u \rangle_2| &= \left| \sum_i \rho_h(i) u(i) \right| \leq \sum_i |\rho_h(i)| |u(i)| \\
&\leq \sum_i |\rho_h(i)| = \|\rho_h\|_1.
\end{aligned}$$

So, if $\sum_h \lambda_h \|\rho_h\|_1^2 < +\infty$, the above inequality and Theorem 10 ensure stability. In addition, since $M_{ij} = \sum_h \lambda_h \rho_h(i) \rho_h(j)$, one has

$$|M_{ij}| \leq \sum_h \lambda_h |\rho_h(i)| |\rho_h(j)| =: f_{ij}.$$

Hence, if $\sum_h \lambda_h \|\rho_h\|_1^2 < +\infty$ one obtains

$$\sum_{ij} |M_{ij}| \leq \sum_{ij} f_{ij} = \sum_h \lambda_h \|\rho_h\|_1^2 < +\infty$$

and this proves also the absolute summability of M .

7.5 Proof of Theorem 12

If there exists $A > 0$ such that $\|\rho_h\|_1 \leq A$ if $\lambda_h > 0$ and the eigenvalues are summable, one has

$$\sum_h \lambda_h \|\rho_h\|_1^2 \leq A^2 \sum_h \lambda_h < +\infty$$

and Theorem 11 then ensures stability.

If the kernel M is stable, its trace is finite and from Lemma 4 we know that the kernel operator \mathcal{M} is compact. Then, as discussed during the proof of Theorem 3, it holds that

$$\text{tr}(M) = \sum_h \lambda_h < +\infty$$

and this concludes the proof.

References

- [1] A. Argyriou and F. Dinuzzo. A unifying view of representer theorems. In *Proceedings of the 31th International Conference on Machine Learning*, volume 32, pages 748–756, 2014.
- [2] N. Aronszajn. Theory of reproducing kernels. *Transactions of the American Mathematical Society*, 68:337–404, 1950.
- [3] K. Atkinson. Convergence rates for approximate eigenvalues of compact integral operators. *SIAM Journal on Numerical Analysis*, 12(2):213–222, 1975.
- [4] C. Baker. *The numerical treatment of integral equations*. Clarendon press, 1977.
- [5] M. Belkin, S. Ma, and S. Mandal. To understand deep learning we need to understand kernel learning. *arXiv e-prints*, Feb 2018.
- [6] S. Bergman. *The Kernel Function and Conformal Mapping*. Mathematical Surveys and Monographs, AMS, 1950.
- [7] M. Bertero, T. Poggio, and V. Torre. Ill-posed problems in early vision. In *Proceedings of the IEEE*, pages 869–889, 1988.
- [8] M. Bisiacco and G. Pillonetto. Kernel absolute summability is only sufficient for RKHS stability. *ArXiv e-prints 1909.02341 2019* <https://arxiv.org/abs/1909.02341>.
- [9] G. Bottegal, A.Y. Aravkin, H. Hjalmarsson, and G. Pillonetto. Robust EM kernel-based methods for linear system identification. *Automatica*, 67:114 – 126, 2016.
- [10] F.P. Carli, A. Chiuso, and G. Pillonetto. Efficient algorithms for large scale linear system identification using stable spline estimators. In *Proceedings of the 16th IFAC Symposium on System Identification (SysId 2012)*, 2012.
- [11] C. Carmeli, E. De Vito, and A. Toigo. Vector valued reproducing kernel Hilbert spaces of integrable functions and Mercer theorem. *Analysis and Applications*, 4:377–408, 2006.
- [12] T. Chen. On kernel design for regularized lti system identification. *Automatica*, 90:109 – 122, 2018.
- [13] 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*, provisionally accepted, 2013.
- [14] T. Chen and L. Ljung. On kernel structures for regularized system identification (ii): a system theory perspective. *IFAC-PapersOnLine*, 48(28):1041 – 1046, 2015. 17th IFAC Symposium on System Identification SYSID 2015.
- [15] T. Chen and L. Ljung. Regularized system identification using orthonormal basis functions. In *2015 European Control Conference (ECC)*, pages 1291–1296, 2015.
- [16] T. Chen, L. Ljung, M. Andersen, A. Chiuso, F.P. Carli, and G. Pillonetto. Sparse multiple kernels for impulse response estimation with majorization minimization algorithms. In *IEEE Conference on Decision and Control*, pages 1500–1505, Hawaii, Dec 2012.
- [17] T. Chen, H. Ohlsson, and L. Ljung. On the estimation of transfer functions, regularizations and Gaussian processes - revisited. *Automatica*, 48(8):1525–1535, 2012.
- [18] T. Chen and G. Pillonetto. On the stability of reproducing kernel Hilbert spaces of discrete-time impulse responses. *Automatica*, 95:529 – 533, 2018.
- [19] A. Chiuso, T. Chen, L. Ljung, and G. Pillonetto. Regularization strategies for nonparametric system identification. In *Proceedings of the 52nd Annual Conference on Decision and Control (CDC)*, 2013.
- [20] Y. Cho and L.K. Saul. Kernel methods for deep learning. In Y. Bengio, D. Schuurmans, J. D. Lafferty, C. K. I. Williams, and A. Culotta, editors, *Advances in Neural Information Processing Systems 22*, pages 342–350. Curran Associates, Inc., 2009.
- [21] F. Cucker and S. Smale. On the mathematical foundations of learning. *Bulletin of the American mathematical society*, 39:1–49, 2001.
- [22] M. Darwish, G. Pillonetto, and R. Tth. Perspectives of orthonormal basis functions based kernels in bayesian system identification. In *2015 54th IEEE Conference on Decision and Control (CDC)*, pages 2713–2718, 2015.
- [23] M.A.H. Darwish, G. Pillonetto, and R. Toth. The quest for the right kernel in bayesian impulse response identification: The use of obfs. *Automatica*, 87:318 – 329, 2018.
- [24] F. Dinuzzo. Kernels for linear time invariant system identification. *SIAM Journal on Control and Optimization*, 53(5):3299–3317, 2015.
- [25] H. Drucker, C.J.C. Burges, L. Kaufman, A. Smola, and V. Vapnik. Support vector regression machines. In *Advances in Neural Information Processing Systems*, 1997.
- [26] N. Dunford and J.T. Schwartz. *Linear operators*. InterScience Publishers, 1963.
- [27] Y. Fujimoto, I. Maruta, and T. Sugie. Extension of first-order stable spline kernel to encode relative degree. *IFAC-PapersOnLine*, 50(1):14016 – 14021, 2017. 20th IFAC World Congress.
- [28] F. Girosi. An equivalence between sparse approximation and support vector machines. Technical report, Cambridge, MA, USA, 1997.
- [29] A. Gittens and M. Mahoney. Revisiting the nyström method for improved large-scale machine learning. *J. Mach. Learn. Res.*, 17(1):3977–4041, 2016.
- [30] P. Heuberger, P. van den Hof, and B. Wahlberg. *Modelling and Identification with Rational Orthogonal Basis Functions*. Springer, 2005.
- [31] G. Kimeldorf and G. Wahba. A correspondence between bayesian estimation on stochastic processes and smoothing by splines. *The Annals of Mathematical Statistics*, 41(2):495–502, 1970.
- [32] S. Kumar, M. Mohri, and A. Talwalkar. Sampling methods for the Nyström method. *J. Mach. Learn. Res.*, 13(1):981–1006, 2012.
- [33] V.B. Lidskii. Non-self-adjoint operators with a trace. *Dokl. Akad. Nauk.*, 1959.
- [34] L. Ljung. *System Identification - Theory for the User*. Prentice-Hall, Upper Saddle River, N.J., 2nd edition, 1999.
- [35] L. Ljung, T. Chen, and B. Mu. A shift in paradigm for system identification. *International Journal of Control*, pages 1–8, 2019.
- [36] H.Q. Minh. Some properties of Gaussian reproducing kernel Hilbert spaces and their implications for function approximation and learning theory. *Constr. Approx.*, 32(2):307–338, 2010.

- [37] B. Ninness, H. Hjalmarsson, and F. Gustafsson. The fundamental role of general orthonormal bases in system identification. *IEEE Transactions on Automatic Control*, 44(7):1384–1406, 1999.
- [38] G. Pillonetto. Consistent identification of Wiener systems: A machine learning viewpoint. *Automatica*, 49(9):2704–2712, September 2013.
- [39] G. Pillonetto. A new kernel-based approach to hybrid system identification. *Automatica*, 70:21 – 31, 2016.
- [40] G. Pillonetto. System identification using kernel-based regularization: New insights on stability and consistency issues. *Automatica*, 93:321–332, 2018.
- [41] G. Pillonetto and B.M. Bell. Bayes and empirical Bayes semi-blind deconvolution using eigenfunctions of a prior covariance. *Automatica*, 43(10):1698–1712, 2007.
- [42] G. Pillonetto, T. Chen, A. Chiuso, G. De Nicolao, and L. Ljung. Regularized linear system identification using atomic, nuclear and kernel-based norms: The role of the stability constraint. *Automatica*, 69:137 – 149, 2016.
- [43] G. Pillonetto and A. Chiuso. Tuning complexity in regularized kernel-based regression and linear system identification: The robustness of the marginal likelihood estimator. *Automatica*, 58:106 – 117, 2015.
- [44] G. Pillonetto, A. Chiuso, and G. De Nicolao. Regularized estimation of sums of exponentials in spaces generated by stable spline kernels. In *Proceedings of the IEEE American Cont. Conf., Baltimore, USA*, 2010.
- [45] G. Pillonetto and G. De Nicolao. A new kernel-based approach for linear system identification. *Automatica*, 46(1):81–93, 2010.
- [46] 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.
- [47] G. Pillonetto, L. Schenato, and D. Varagnolo. Distributed multi-agent Gaussian regression via finite-dimensional approximations. *IEEE Trans. on Pattern Analysis and Machine Intelligence*, 41(9):2098–2111, 2019.
- [48] T. Poggio and F. Girosi. Networks for approximation and learning. In *Proceedings of the IEEE*, volume 78, pages 1481–1497, 1990.
- [49] D. Robert. On the traces of operators (from Grothendieck to Lidskii). *EMS newsletter*, 2017.
- [50] W. Rudin. *Real and Complex Analysis*. McGraw-Hill, Singapore, 1987.
- [51] B. Schölkopf, R. Herbrich, and A. J. Smola. A generalized representer theorem. *Neural Networks and Computational Learning Theory*, 81:416–426, 2001.
- [52] Hongwei Sun. Mercer theorem for RKHS on noncompact sets. *J. Complexity*, 21(3):337–349, 2005.
- [53] V. Vapnik. *Statistical Learning Theory*. Wiley, New York, NY, USA, 1998.
- [54] G Wahba. *Spline Models For Observational Data*. SIAM, Philadelphia, 1990.
- [55] B. Wahlberg. System identification using Laguerre models. *IEEE Transactions on Automatic Control*, 36(5):551–562, 1991.
- [56] B. Wahlberg. Laguerre and Kautz models. *IFAC Proceedings Volumes*, 27(8):965 – 976, 1994. IFAC Symposium on System Identification (SYSID’94), Copenhagen, Denmark, 4-6 July.
- [57] C.K. I. Williams and M. Seeger. Using the Nyström method to speed up kernel machines. In *Proceedings of the 2000 conference on Advances in neural information processing systems*, page 682688, Cambridge, MA, USA, 2000. MIT Press.
- [58] E. Zeidler. *Applied Functional Analysis*. Springer, 1995.
- [59] H. Zhu, C.K.I. Williams, R.J. Rohwer, and M. Morciniec. Gaussian regression and optimal finite dimensional linear models. In C. M. Bishop, editor, *Neural Networks and Machine Learning*. Springer-Verlag, Berlin, 1998.